

# DAVIDSON'S TEST

## Donald Davidson's Critique of the Turing Test as an Expression of his Theory of Intellectual and Linguistic Competence

TOBIAS LOHSE

NOVEMBER 22, 2015

### Abstract

In this paper I discuss Donald Davidson's critique ([1990b] 2004) of the Turing Test in which he agrees with Turing (1950) that communication is the essential hallmark of human-like intellectual abilities and proposes a modified version of the Turing Test that requires the judge to observe the learning process of the computer through social interactions in a shared world. This modified test serves as an illustration of Davidsons ideas about the nature of our linguistic and intellectual competences, especially his Epistemology of 'Triangulation' of subjective, intersubjective, and objective knowledge ([1991b] 2001) and his empirical 'Unified Theory of Meaning, and Action' ([1980] 2004). My discussion aims to expose how we can treat Davidson's Test as an operational definition of his concept of intellectual competence and how it relates to his rejection of other theories of meaning.

Furthermore, I argue that Davidson and Turing agree that the essence of communication can be captured in a formal theory and that both saw social interaction in a shared world as essential for human-like intellectual competence; for Davidson this is grounded in his idea of a triangular Epistemology; for Turing it is grounded in his understanding that Machine Learning is needed to achieve Artificial Intelligence. Davidson however goes further by denying that such experience of interaction can simply be encoded in predetermined conventions, and requiring the ability for 'Radical Interpretation' of previously unknown expressions in the very instance of communication. In the end, I give an outlook for the implications of Davidson's thoughts for Machine Learning and Computer Linguistic approaches that would deserve further research.

### Keywords

Donald Davidson, Philosophy of Mind, Philosophy of Language, Epistemology, Interpretation, Turing Test, Artificial Intelligence, Computer Linguistic

# Contents

<b>Introduction</b>	<b>1</b>
Introduction to Donald Davidson . . . . .	1
<b>1 The Turing Test and Its Relevance</b>	<b>1</b>
1.1 Interpretations and Versions of the Turing Test . . . . .	2
1.2 Relevance of the Test and the Turing Principle . . . . .	2
1.3 Turing on Machine Learning . . . . .	3
<b>2 Davidson's Critique of the Turing Test</b>	<b>4</b>
2.1 Relevance of Turing's Test for Davidson's Philosophy of Mind . . . . .	4
2.2 Davidson's Critique . . . . .	5
2.3 Davidson's Proposal for a Modified Test . . . . .	6
<b>3 Davidson's Theory of Intellectual and Linguistic Competence</b>	<b>7</b>
3.1 The Meaning of Programs and Anomalous Monism . . . . .	7
3.2 Triangulating Knowledge of a Shared World . . . . .	9
3.3 From Prior to Passing Theories of Interpretation . . . . .	11
3.4 Davidson's Test and the Theory of Intellectual and Linguistic Competence . . . . .	14
<b>Conclusion</b>	<b>15</b>
Davidson's Theory and Machine Learning . . . . .	15
Implications for Artificial Intelligence and Computer Linguistics . . . . .	15
<b>Notes</b>	<b>15</b>
<b>Acknowledgements</b>	<b>19</b>
<b>References</b>	<b>19</b>

# Introduction

This paper discusses Donald Davidson's critique ([1990b] 2004) of the Turing Test as an illustration of Davidson's theory of linguistic and intellectual competence. After a short Introduction to Donald Davidson it consists of three main sections.

**Section 1** explores Turing's original ideas and their relevance. It discusses different interpretations and versions of the Turing Test and establishes which interpretation I follow. It also contains a discussion of the significance of the test and Turing's ideas about Machine Learning for the field of Artificial (General) Intelligence.

**Section 2** presents and evaluates where Davidson agrees and disagrees with Turing. It also presents Davidson's modified version of the Turing Test which allows to judge the computer based on an observation of its learning process. I argue that this test acts as a definition for intelligence and expression of Davidsons theory of linguistic and intellectual competence.

**Section 3** delves deeper into Davidson's theory by looking at its parts: reconstructing his argument for supervenience of the mental, exposing his epistemology of triangular knowledge, and describing his empirical anti-conventionalist theory of interpretation. It shows how those elements come together as a unified theory about the essence of linguistic and intellectual abilities that constitute Davidson's definition of rationality or intelligence that is incorporated in Davidson's Test. After that I discuss how Davidson's Test fences against some classical objections against the Turing Test and some other concerns.

I conclude with a summary of the core agreements and distinctions between Davidson and Turing and a brief discussion of the implications of Davidson's ideas (holistic theory and anti-conventionalism) for Computer Linguistic and Machine Learning approaches that might deserve further investigation.

## Introduction to Donald Davidson

Donald Davidson (1917–2003) “was one of the most important philosophers of [action, language, and mind in] the latter half of the twentieth century [...] with a reception and influence that, of American philosophers, is perhaps matched only by that of [his friend, mentor, and colleague] W. V. O. Quine” (Malpas [1996] 2015, p. 1). His thoughts are exposed through a collection of very densely written essays (published from 1963–2001)<sup>1</sup> which “form a mosaic out of which emerges a unified and surprisingly elegant overall view of the mind and its relation to the world. It sees our *nature as linguistic beings* as the key to the possibility of thought, to the objectivity of the world [...], and to how the [irreducibly mental] mind moves us to action in a [physical] world” (Lepore and Ludwig 2009, p. 1, emphasis added).

To give a very brief overview of Davidson's philosophy:



Figure 1 – Donald Davidson by David Levine in Forum Gallery, New York

- He takes the causal nature of thought ([1963] 2001) as an argument for the supervenience and irreducibility of the mental to the physical known as ‘Anomalous Monism’ ([1970] 2001).
- He argues that semantic theories for natural languages must be recursive ([1965] 2001) and that meaning and propositions can be replace by truth conditions and distal stimuli ([1967] 2001).
- Based on Quine's Radical Translation, he takes interpretation as a process of adapting a theory of intertwined understanding and believe about the speaker ([1973b] 2001).
- He rejects Conventionalism ([1984] 2001) based on his argument for the primacy of ideolect, placing the social aspect of language in Epistemology ([1994] 2005).
- Rejecting positivist and cartesian Epistemology, he sees knowledge of self, world, and others as interdependent ([1991b] 2001).
- His Unified Theory ([1980] 2004) aims to captures the essence ([1995a] 2004) of mind, language, and action.
- He discusses similarities between his philosophy and that of Spinoza ([1999] 2005) and Plato and Gadamer ([1997a] 2005).

## HISTORICAL CONTEXT

Language has been a philosophical topic since Plato, but it especially gained attention in the so-called ‘Linguistic Turn’ of the 20th century, when it was recognized as the medium for our rational access to the world. Names like Frege, Russel, the early Wittgenstein, and the Vienna Circle stand for the earlier Ideal Language Philosophy and Logical Empiricism which hoped to solve problems of Metaphysics through a *definite logical analysis* of meaning. However, This approach did not turn out very successful. Led by Austin, Ryle, and the late Wittgenstein analytical philosophy paid more attention to Ordinary Language, focusing on the importance of *social and performative aspects*. Quine and Davidson stand at the resolution of this tension with an approach that recognizes linguistic behavior in its context while leveraging the formal methods for their analysis (Bertram 2011). Through Davidson the analytical tradition also became linkable to the ‘continental’ Hermeneutic traditions, contributing to the dissolution of the philosophical divide of the 20th century (Malpas [1996] 2015, p. 2).

## 1 The Turing Test and Its Relevance

Alan Turing (1912–1954) first proposed his famous test in ‘Computing Machinery and Intelligence’ (1950). The thesis of his paper is that

the question “Can machines think?” — which he deemed too vague to deserve discussion — can be replaced with the question whether a computer can pass a specified test (that is now known as the ‘Turing Test’). This test is described as a game: the imitation game.

The imitation game consists of an *interrogator* who can communicate via a real-time text chat interface with two players. One player is a *computer*, the other a *human*. The interrogator’s task is to identify the human after a given time. The computer’s task is to pretend to be a human<sup>2</sup> and trick the interrogator into wrongly identifying it as the human. The human’s task is to help the interrogator to make the correct identification. If many interrogators consistently<sup>3</sup> can’t tell the computer apart from the human the computer wins the game and passes the test. (Turing 1950, p. 433–434)

In the following I discuss different interpretations of the Turing Test and establish which interpretation and version of the test I take as the basis for my discussion. Furthermore, I argue for the relevance of this interpretation of the test and work out 4 main claims from Turing’s discussion of his test. In doing so I specifically emphasize some of Turing’s visionary ideas on how a computer must be programmed to pass the test, that have implications for Davidson’s discussion of the test.

## 1.1 Interpretations and Versions of the Turing Test

The interpretations are not as straight forward, as it might seem. There are three main interpretations for what it means to pass the Turing Test:

1. The Turing Test is taken as an operational *definition* — a *sufficient* and *necessary* condition — of *intelligence*<sup>4</sup>: Something is intelligent if and only if it passes the Turing Test. (Such an interpretation is for example found in Millar 1973.)
2. The Turing Test is taken as a *sufficient* condition of intelligence: Something is intelligent if it passes the Turing Test. But it is not necessary for something to pass the Turing test in order to be intelligent. (Davidson [1990b] 2004 takes this interpretation.)
3. The Turing Test is taken “as a potential source of good inductive evidence for” intelligence: If Something passed the Turing Test one would be *justified* for inductively inferring that it is intelligent. (This interpretation goes back to Moor 1976, p. 249, 251.)

The first interpretation and any reading of intelligence in a broader sense are clearly not in line with Turing’s ideas. Turing was not interested in a definition of intelligence but in setting a clear goal for further research. In fact I would argue that the third interpretation from Moor does Turing most justice. This becomes fairly clear from the following section of a radio interview ([1952] 1999), which represents Turing’s simplest expression of the test:

“I *don’t want to give a definition of thinking* [...] I don’t really see that we need to agree on a definition at all. The important thing is to try to *draw a line between the properties of a brain, or of a man, that we want to discuss, and those that we don’t*. [...] I would like to suggest a particular kind of test that one might apply to a machine. You might call it a test to see whether the machine thinks, but it would be better to avoid begging the question, and *say that the machines that pass are (let’s say) ‘Grade A’ machines*. The idea of the test is that the machine has to try and pretend to be a man, by answering questions put to it, and it will only pass if the pretence

is reasonably convincing.” (Turing [1952] 1999, p. 466, emphasis added)

Turing clearly describes the imitation game as a test that defines a special *class of computers*, but not as one that defines thinking/intelligence. It is also implicit that he sees the *ability to communicate* as quintessential for human-like intellectual abilities.

The quotation also gives a simplified version of the test, in so far as it removes the second human player. While this might be limiting for a quantitative analysis, it ensures the focus on the main question and removes any undue emphasis on strategy. Davidson proposes the same simplification in his discussion ([1990b] 2004). He goes even further in proposing that since “we might count an object as thinking even if it were easily distinguishable from a person. [...]he interrogator [...]should simply] be asked to decide whether or not the object is thinking” ([1990b] 2004, p. 81). If we assume that any bias against the object’s ability to think can be removed (which makes sense for this philosophical discussion), this seems to be the clearest expression of Turing’s idea. Therefore, I will follow Davidson and refer to this simplified version in the following.

## 1.2 Relevance of the Test and the Turing Principle

Where does it leave the relevance of the Turing Test, if we follow Moor and interpret it merely as a framework to gather evidence for human-like intelligence? For Computer Scientists it might be *prudent* to follow Turing and ignore the philosophical question of a definition of intelligence and instead take the Turing Test as a definition of a certain class of computers. In this way the Turing Test can give a clear goal to the field of ‘Artificial General Intelligence’ (AGI), for which the unclarity about a clear definition of its topic is a core challenge. Unlike other propositions for an operational definition<sup>5</sup> of AGI, the Turing Tests provides a clear and well-justified empirical goal. In his analysis Moor argues why the Turing Test is well apt for that:

“There are two strong arguments why the Turing test is a good format for gathering inductive evidence. First, the Turing test permits direct or indirect testing of virtually all of the activities one would count as evidence for thinking. Secondly, the Turing test encourages severe testing. [...] The computer would be tested in detail over a wide range of subjects [...]and] the interrogator’s goal is to find a refuting instance which gives the computer away.” (Moor 1976, p. 251–252)

Critics of the Turing Test as a goal for AGI mostly fall into two camps. The first point Moor provides is aimed against the first type of critic who suggests the test sets the wrong goal. As Moor argues, communication is a very clear framework to investigate all kinds of thinking. A further practical reason for the relevance of the Turing Test — that I would add — is that natural language communication provides a ‘gold standard’ for completely natural and seamless computer user interfaces. Many of those critics mistake the question as philosophical while it is best treated as definitional. Computer Scientists ought not to be concerned with defining intelligence in general but with a good definition for “Grade A” computers — to use Turing’s terminology. In this regard, the interpretation of the Turing Test as a framework to gather empirical evidence can be set as a clear definition of the goal of AGI.

The other type of critic questions the adequacy of the test to determine whether the goal is reached. Moor’s second point is aimed

against that by pointing out how well the test encourages thorough testing. The critics sometimes mistake Turing's predictions as a specification that the test ought to take only 5 minutes and limited implementations of the Turing Test that favor engineering tricks like the Loebner Prize<sup>6</sup> discredit Turing's intentions. I agree with Copeland's interpretation of Turing: the goal set by the test is a computer that "plays the imitation game successfully come what may, with no field of human endeavour barred, and for any length of time commensurate with the human lifespan." (Copeland 2000, p. 530)

### TURING'S THREE MAIN CLAIMS

Following Copeland (2000, p. 530), we can establish the following two main claims of Turing from our previous discussion:

1. Turing sees communication as suitable to expose the relevant intellectual abilities that determine whether a computer can perform tasks on par with a human being.
2. Turing thinks that his test *specifically* can determine whether a computer possesses such communicative abilities.

But there is also a third claim which is known as the Turing Principle<sup>7</sup>:

3. Turing believes that universal computers can simulate any physical process, including the brain.

This might be most clearly expressed in Turing's lecture on 'Can Digital Computers Think?': "If it is accepted that real brains [...] are a sort of machine it will follow that our digital computer suitably programmed will behave like a brain" (Turing [1951] 1999, p. 463). This is a much more controversial claim and — as Turing was well aware — he had few arguments and even less evidence for it.

It is important to highlight that this claim is not central in Turing's writing and that it is separate and fully *independent* from the others claims. Nevertheless, it is a point philosophers have attacked. (The most famous example probably being Searle 1980.) This is often tied to the misinterpretation of the Turing Test as a definition or sufficient condition for intelligence in a philosophical sense. I am not particularly interested in the discussion of this third claim here.

### 1.3 Turing on Machine Learning

Unlike the obsession in popular culture with human-like 'Artificial Intelligence' (AI) might suggest, the interest of Computer Scientists in the last years has been more focused on so-called 'weak' AI. This refers to applying 'Machine Learning' to apply domain-specific intelligence to special problems (most notably: image recognition, domain-specific language processing, and robotics). The field that is concerned which this is referred to as 'Artificial Intelligence' these days. The field that is concerned with human-like strong AI is the field of Artificial General Intelligence (AGI) which was referenced in the previous section. This shift of focus has mainly been due to a lack of success in AGI.

Turing believed that by the turn of the century there would be computers that "play the imitation game so well that an average interrogator will not have more than 70 per cent chance of making the right identification after five minutes of questioning" (1950, p. 442).

Recent winners of the Loebner Prize (AISB 2015), which is awarded to the most human-like chatbot judged by a Turing-inspired test, show that this has not quite come to pass. On the other hand, Turing himself said that it would take "*at least* 100 years" ([1952] 1999, p. 434) until a general Turing Test could be passed — and the goal of AGI reached — which is still well within the realm of the possible. Turing's wrong prediction of the development can be mainly accredited to his underestimation of the required processing power for Machine Learning (1950, p. 455). Recent accomplishments in AI with deep neural networks, for example, have only become feasible because computers are able to run networks with millions of neurons in real time.<sup>8</sup>

### TURING'S FOURTH MAIN CLAIM

However, Turing's idea about *how* to build a computer that could pass his test are more interesting than his timeline predictions. His claim that "the problem [of building a computer that passes the Turing Test] is mainly one of programming" (1950, p. 455) still rings true today.

We have already learned of his idea that computers should be able to simulate the brain ([1951] 1999). Indeed this reverse engineering approach to AI is the basic idea behind the neural networks which power the most successful image entity-recognition algorithms today. Turing ([1948] 1992) pioneered this approach with his B-type unorganized machine which consisted of neurons that are trained through an 'education' process, and he proved that they were equivalent with digital computers (see also Copeland and Proudfoot 1999). It is plausible that we might achieve AGI through brain simulation before we even deeply understand how the brain works.

Even more interesting is that Turing also predicted the Machine Learning approach to build intelligent algorithms and saw its non-deterministic nature as a characteristic feature:

"Instead of trying to produce a programme to simulate the adult mind, why not rather try to produce one which simulates the child's? If this were then subjected to an appropriate course of education one would obtain the adult brain. [...] We have thus divided our problem into two parts. The child-programme and the education process. These two remain very closely connected. We cannot expect to find a good child-machine at the first attempt. One must experiment with teaching one such machine and see how well it learns." (Turing 1950, p. 456)

"An important feature of a learning machine is that its teacher will often be very largely ignorant of quite what is going on inside, although he may still be able to some extent to predict his pupil's behaviour. [...] This is in clear contrast with normal procedure when using a machine to do computations: one's object is then to have a clear mental picture of the state of the machine at each moment in the computation. [...] Processes that are learnt do not produce a hundred per cent certainty of result; if they did they could not be unlearned." (Turing 1950, p. 458-459)

This description is very much how modern statistical Machine Learning algorithms work. They consist of an algorithm that describes a mathematical model which is trained with human-annotated data (for example, a grammar analysis of sentences) and are then able to perform their task on similar data. (See Schubert [2014] 2015 for an overview about approaches to Natural Language Processing; and Jurafsky and Martin 2015 for a more technical introduction.) As

Turing mentions, this is a paradigm shift from classical algorithms which have results that are clearly defined by the programmer and predictable independent of any training data. We might add a fourth core claim:

4. Turing believes that the approach to device intelligent computers should be based on learning algorithms that are trained and do not behave predictable in a classic sense.

Note that devising intelligent algorithms is a quite different task and involves techniques where an interpretation of the state of the program at each step becomes difficult or even impossible — especially in the case of deep neural networks. And where the outcome is not just dependent on the set of rules specified in the programming but also on the ‘experience’ gathered in the program during its learning process. We will see later why this matters for a Davidsonian perspective.

## SUMMARY

We have seen that a simplified interpretation of the Turing Test as proposed by Davidson ([1990b] 2004) reveals the core of Turing’s idea. However, we have also seen that an interpretation of the Turing Test as an operational definition or sufficient condition for intelligence is not in line with Turing’s writing. Instead, I have proposed to follow Moor (1976) and interpret the test as a framework to collect empirical evidence to show that a computer can perform human-like tasks. We have seen that this can provide a good and clear definition of a special class of computers that are the goal of AGI research.

Furthermore, we have established that Turing claims that (1) communication abilities are representative of intellectual abilities in general and (2) that his test is adequate to evaluate those abilities. I have pointed out that his claim (3) that computers can simulate the brain (the Turing Principle) is completely independent of the other claims and not very essential for Turing.

Lastly, we have seen that Turing had pioneering ideas about how to devise algorithms which could pass his test. Such algorithms (4) need to be able to learn when trained with data and, different from classical algorithms, their outcome is not predictable independent of their ‘experience’ and the states of their operation are not easily interpretable.

## 2 Davidson’s Critique of the Turing Test

In this section I will first discuss the relevance of the Turing Test for Philosophy of Mind and for Davidson specifically. I will then argue where Davidson agrees with Turing, where his criticism of Turing applies, and in how far I think this criticism does Turing justice. The section ends in an exposition of Davidson’s proposal for a modified Turing Test that provides a definition for intelligence.

### 2.1 Relevance of Turing’s Test for Davidson’s Philosophy of Mind

I have argued that the Turing Test is misunderstood as a definition or sufficient condition for intelligence in general and merely provides a pragmatic goal for AGI research. But why should it be relevant

to Philosophy of Mind then? Davidson says “the test is designed to throw light on the nature of thought [...and it] can be applied to any object” ([1990b] 2004, p. 78). While Davidson interprets the test as a sufficient condition for thought, his statement also reveals another more important aspect that is echoed by Turing:

“The whole thinking process is still rather mysterious to us, but I believe that the attempt to make a thinking machine will help us greatly in finding out how we think ourselves.” (Turing [1951] 1999, p. 465)

So the Turing Test and the related task of constructing a machine which can pass it, is a practical device to learn more about the nature of thought and communication. This resonates very much with Davidson’s philosophy. For the essential paradigm shift of his philosophy about language is following Quine’s rejection of Carnap’s Logical Empiricism<sup>9</sup>. The objective of a theory of meaning for Davidson “is sought not in reductive analyses [of the meaning of expressions], but rather in showing how evidence can be marshalled in support of a theory<sup>10</sup> of interpretation for a speaker.” (Lepore and Ludwig 2009, p. 9) Philosophy is not an a priori discipline but simply a more general empirical discipline. “‘Language is a social art’ [...and] evidence for its acquisition and deployment must be intersubjective, and, hence, recoverable from overt behaviour” (Lepore and Ludwig 2009, p. 22–23). Therefore the search for a definition of meaning or of intelligence itself is rather uninteresting. A definition is only relevant as part of an empirical theory that can be judged by its successful application. (I’ll talk more about this at the end of this section.) The search for a definition is therefore replaced by the search of an empirical theory of intelligence.

From this angle the task of programming a machine which can pass the Turing Test, seems closely related to that of devising and testing an empirical theory of intelligence. At least if we accept that the essence of such a theory could be expressed recursively. Judging a computer as passing the test is not so different from judging a person as rational being after all. Moor argues:

“I believe that another human being thinks because his ability to think is part of a theory I have to explain his actions. [...] The evidence for the theory comes from the outward behavior of the person. [...] There is no reason why knowledge of computer thinking can not arise in the same way. I can use the computer’s behavior as evidence in assessing my theory about its information processing.” (Moor 1976, p. 251)

Davidson and Turing agree that we can devise a scientific theory that describes the essential parts of our linguistic competence. In the case of Turing this is evident in his belief that there is a program that allows the computer to win the imitation game. For Davidson it is evident in the proposal of his Unified Theory which can capture the essence of linguistic competence and rationality. However there are some important distinctions between the nature of these approaches that will be highlighted in subsection 3.1.

### LINGUISTIC ABILITIES AS THE ESSENCE OF INTELLIGENCE

What we have discussed so far depends on accepting Turing’s first claim that communication abilities are representative of intellectual abilities. As mentioned in the introduction, Davidson agrees since he “sees our nature as linguistic beings as the key to the possibility

of thought" (Lepore and Ludwig 2009, p. 1). This belief springs from his metaphysics<sup>11</sup> and from his epistemology. The latter will be examined in more detail in subsection 3.2. However, Davidson emphasizes the interrogator's judgement as the essential point:

"Turing was right, in my opinion, in taking as the only test for the presence of thought and meaning the interpretive powers and abilities of a human interpreter." (Davidson [1990b] 2004, p. 86)

As we will discuss in subsection 2.3 Davidson very intentionally changes the emphasis here, since he argues that the ability to be interpreted is not a sufficient, but a necessary condition for attributing thought — introspection into the working of the mind for example is *not* a sufficient condition for Davidson. For Davidson the essence of Turing's approach is that "instead of asking how the content of a concept [...] is thought of by the creature that has the concept [...], we ask [...] how an observer can size up the contents of the thoughts of another creature" (Davidson [2001] 2004, p. 137).

The third claim (Turing's Principle) is not of particular interest to Davidson, but because of his naturalistic ontology<sup>12</sup> he agrees: "A person is a physical object which [...] functions according to physical laws. So [...] there is no reason why an artificial object could not think[...] . The real question is:] how much like us must an artifact be, and in what ways, to qualify as having thoughts?" ([1990a] 2004, p. 87).

The second claim (the appropriateness of the test) is what Davidson discusses and criticizes extensively in his essay 'Turing's Test' ([1990b] 2004). The fourth claim (the importance of Machine Learning) is something Davidson does not pay much attention to, but that is closely related to his critique. Both are the topic of the next subsection.

So we have seen why *Davidson* finds Turing's Test particularly interesting. I have chosen his critique for discussion here, since it clearly relates to Computer Science and because it "opens the way for Davidson's own view into the nature of thought" (Cavell 2004, p. xvii).

## 2.2 Davidson's Critique

Davidson finds Turing's Test inadequate to show that an object is thinking. Not because communication is not a sufficient criterion for intellectual competence for Davidson — as we have discussed. Neither because a test would require an introspection into the workings of the mind. Davidson explicitly says that it is not inadequate because it "restricts the available evidence to what can be observed from the outside" ([1990b] 2004, p. 83). But because it does not enable the interrogator to observe a history of three-way engagement between the object, a shared world, and other minds in which the object develops its semantics.

This is related to the "fundamental difference between semantics, which relates words to the world, and syntax, which does not" (Davidson [1990a] 2004, p. 94, he adopts the terminology from Tarski). For Davidson, the essential requirement for thought or intelligence is to assign meaning to words, to relate them to the world. But he argues that the interrogator in the Turing Test can not guarantee that the computer is able to do that:

"The interrogator [...] has no clue to the semantics of the object. There is no way he can determine the connection between the

words that appear on the object's screen and events and things in the world. Of course there must be some connection; there is no other way to account for the intelligibility of the object's English. [...But] it is perfectly possible that the connection between words and things was established by someone who programmed the object, and then provided purely syntactic connections between words for the object to wield. In this case it is the programmer who [...] has given meaning to the words [...], but the object doesn't mean anything, and there is no reason to take it to be thinking.

In order to discover whether the object has any semantics, the interrogator must learn more about the connections between the output of the object and the world [...], through observing relevant causal interactions between the speaker, the world, and the speaker's audience. [...] The interrogator [must be allowed] to watch the object interact with the world." (Davidson [1990b] 2004, p. 83)

The interpreter can only determine whether the computer means something by its words if he or she can tell what the computer means by them. Therefore "any evidence that thinking is going on will have to be evidence that particular thoughts are present" ([1990b] 2004, p. 80). While it is clear that thoughts have caused the computer's interaction, the test framework is not giving the interrogator a chance to determine how any thought has come to the computer's knowledge, only that it has.

Davidson accredits this failure to the fact that "Turing wanted his test to draw '... a fairly sharp line between the physical and the intellectual capacities of man' (p. 434). [But t]here is no such line" (Davidson [1990b] 2004, p. 84). Turing was wrong that the 'body' of the object does not matter. Even though the details in which the sense organs convey impressions may not matter, the existence of such organs can not be reduced to purely textual communication. Turing has a somewhat Cartesian<sup>13</sup> approach to thought. He might not prescribe to an ontological mind-body dualism, but he tends to view thoughts as independent from sensory access to the world and only dependent on being able to communicate. Turing claims the "example of [the deaf blind] Helen Keller shows that education can take place provided that communication in both directions between teacher and pupil can take place" (Turing 1950, p. 456). He overlooks that Helen Keller for one had been able to see and hear for her first 19 months and, more importantly, only learned to communicate when her teacher correlated her signs to Helen's feeling of touch (Wikipedia 2015). So this example might serve more to illustrate Davidson's externalist approach to epistemology rather than underline Turing's point that communication alone is essential for thought.

The ability to interact with the world in any way and a history of such interactions and communications about those experiences with others, is essential for developing a semantics. You "don't understand a language if there are not numerous connections between your use of words and experiences" (Davidson [1990b] 2004, p. 85). The interrogator needs to observe the computer's history of engagement with the world, to judge its intelligence.

### FAIRNESS OF THE CRITIQUE

So Davidson's main critique is that Turing's idea of a sharp line between physical and intellectual abilities is ill-conceived and Turing doesn't see the importance of interaction in a shared world as essential for (1) *developing* thoughts and for (2) *judging* intellectual abilities.

While the second critique is certainly justified, since the computer is hidden from the interrogator, the first point is debatable in light of Turing's work about Machine Learning. In his last footnote, Davidson claims that Turing "views this [learning] simply as an economical way of producing a device with mature thoughts; he does not see it as the only way" ([1990b] 2004, p. 86). As outlined in section 1, I think this is overestimating Turing's interest in presenting a sufficient condition for intelligence and underestimating his interest in actually building a certain grade of computer. Turing's ideas on Machine Learning take up a third of his essay (1950, p. 454–460) and constitute — as I argued — a fourth main claim. Most of his research from 1948–1952 was indeed not focused on the test which became so famous, but on approaches to build learning computers and simulating neural networks. None of this work mentions any other approaches for building AI. Turing actually outlines ideas how a computer with an indexed memory can be constructed that relates past experiences to new situations using associative connections and evaluation based on reactions of its teacher and later based on self constructed norms (see Turing [1951] 1996, p. 257–258). This seems closer to Davidson's requirement of having a history of interaction in a shared world than he gives credit for. So maybe classifying Turing as having a Cartesian approach to epistemology is premature.

Nevertheless, Davidson's critique of the *test* is justified. In the light of our discussion of Turing's work however, our perspective might change: from seeing the inadequacy of the test as a result of Turing's underestimation of the importance of learning and interaction with the world; to surprise that he deemed the learning process, that he saw as essential too, as so irrelevant for the judgement.

## 2.3 Davidson's Proposal for a Modified Test

Davidson in fact, does not simply critique the test, but he proposes a modified test:

"The Test must be modified [...]. The object must be brought into the open so that its causal connections with the rest of the world as well as with the interrogator can be observed by the interrogator.

Can the interrogator now tell what the object thinks? The answer is that it depends [...]. Let us suppose the interrogator finds that the object uses words just as he does [...and] infers (let us suppose correctly) that the object's linguistic dispositions are similar to his own in relevant ways. In the case of a person, the interrogator would be justified in assuming that these dispositions were acquired in the usual way: in the basic cases, by past causal intercourse with things and circumstances of the sort to which the person is now disposed to respond. [...] But the assumption is not justified in the case of a computer: [...] The computer which has never experienced a dog and has no memory of dogs can't mean dog by the word 'dog' [...]. Thought and meaning require a history of a particular sort. [...] Unless we [...] can observe it in action over time, we have no basis for guessing how a computer came to have the dispositions it has.

It is unclear exactly what kind of history is necessary [...]. But our intuitions are clear enough in many cases. You [...] don't understand a language if there are not numerous connections between your use of words and experiences[. ...] It may seem that minds are, after all, inscrutable if no present observation of their operation can reveal what they are thinking. But of course this does not follow. [...] Even the mind of an artefact can, if it has one, be understood; it just takes longer, long enough for some history to be observed, since it cannot be inferred." (Davidson [1990b] 2004, p. 84–86)

Davidson's rejection of the "sharp line between the physical and the intellectual capacities" ([1990b] 2004, p. 84) should not be misinterpreted as a rejection of the possibility to devise a test at all. In his follow up essay 'Representation and Interpretation' ([1990a] 2004) he comes back to the question "what could we detach from a person and still count him or her as a thinking creature" (p. 87). And agrees with Turing that origin, building material, and size and shape for example are irrelevant. In so far, Turing was on the right track with his test, he only went too far in detaching the history from a person. The ability to determine non predefined meaning — to connect symbols with objects and events in the world — is essential for intelligence and it depends on a rich base of experiences of causal interactions with these events and objects.

I put Davidson's explanation of the required nature of the test as follows: We may judge an object as intelligent, if (and only if — as I will argue) the object can be observed to be able to

1. successfully communicate with other intelligent beings and
2. derive its own semantics from a history of its experiences of interaction with other intelligent beings and with objects and events in a shared world.

This is what I call 'Davidson's Test'.

The ability to come up with its own semantics from a history of interaction ensures that the computer has a rich conceptual system in Davidson's holistic approach. Because "to have even one thought — one belief or desire — a computer would have to have a very great many other thoughts and desires. Beliefs and desires can exist only in the context of a very rich conceptual system." (Davidson [1990a] 2004, p. 90)<sup>14</sup>.

We will see in how far this test is an expression of Davidson's theory of linguistic and intellectual competence throughout the next section and specifically in subsection 3.4.

## INTERPRETATION OF DAVIDSON'S TEST AS A DEFINITION

So we have a definition of Davidson's Test, it remains to explain how it is to be interpreted. Davidson takes the Turing Test as aiming "to discover whether a sufficient condition for thought is satisfied; the condition is not claimed to be necessary" (Davidson [1990b] 2004, p. 81). So we can certainly interpret Davidson's Test as a sufficient condition for intelligence. But Davidson seems to go even further:

"The *only* way to tell if an artificial device [...] has [thoughts...] and the ability to perceive and interact with the world as a person does, is to attempt to *interpret the behavior* of the device in the same way we do the behavior of a person. [...] Understanding the program and physics of a device [...on the other hand] is *not* [...] sufficient for] understanding the thought [...] of that device." (Davidson [1990a] 2004, p. 99, emphasis added)

It seems to me, that this is interpreted correctly as the claim that Davidson's Test provides not only a sufficient, but also a necessary condition — and therefore constitutes an operational definition of intelligence. Interpreting the behavior of a person is nothing else for Davidson than interpreting its utterances (or some other communications with semantic content). This follows from the fact that the interpretation of the semantic content of an expression requires the interpreter to have a theory about the beliefs and semantics of the person and this theory must be based on the behavioral evidence the



interpreter has. We will investigate this theory of interpretation from Davidson in more detail in subsection 3.3.

The second part of Davidson's claim aims to reflect an obvious counter argument that might occur. If one completely understood the inner workings of a device's or person's brain, one could certainly also judge whether it was intelligent. While Davidson does not want to deny that dissecting a person's brain might enable one to undoubtedly judge a person as belonging to the species *homo sapiens* and therefore be justified to infer that he or she is intelligent; this is not the kind of proof we are looking for. The test provides an *explicit* definition of intelligence, while there might be other *implicit* ways of inference about intelligence. Davidson *does* reject that any knowledge of the brain can *explicitly* prove intelligence. To illustrate this he takes the programming of a machine as an example, which I will discuss in the next subsection.

I have argued before that the search for a definition is replaced by the search for an empirical theory in Davidson's philosophy. So if we talk about an operational definition here, this only becomes relevant holistically because it is part of an empirical theory of intelligence — and in Davidson's case this means a theory of interpretation.<sup>15</sup> So this theory is expected to fit empirical data about interpretations, to allow us to apply it and actually yield correct interpretations of utterances, or (since the theory is not developed far enough yet for that) at least to fit our intuitions for how we interpret others. Therefore we can't judge this definition by itself; we need to look at Davidson's theory as a whole and see whether it is a convincing empirical theory for our linguistic and intellectual competence, then we can judge whether this definition of intelligence, which provides the goal for this theory, is convincing to us. Hence we will discuss Davidson's theory in more detail in the next section.

#### SUMMARY

We have started with a discussion of how Davidson's approach to language as a social art and interpretation as an empirical process relates to Turing's approach to intelligence; pointing out that both Davidson and Turing think that intelligence can be described empirically and that there is a relation between the process of building an intelligent machine and devising an empirical theory of intelligence. I have also argued that both see linguistic abilities as essential for intelligence and agree in their affirmative answer to the question whether artificial intelligence is possible, even though both do not find the question very interesting.

Then we have learned that Davidson critiques Turing's Test for not exposing how the computer develops its own semantics from a history of experiences of interactions in a shared world. I have argued that this shortcoming of Turing's design is surprising because he agrees with Davidson that learning from experience is essential for developing intelligence.

At the end I have given a characterization of Davidson's Test that requires not only the observation (1) of successful communication but also (2) of the process of deriving semantics from a history of interactions in a shared world. I argued that this test should be interpreted as an operational definition of intelligence which we can only judge in light of how empirically convincing Davidson's theory of linguistic and intellectual competence is for which it provides the goal.

## 3 Davidson's Theory of Intellectual and Linguistic Competence

The goal of this section is to relate Davidson's arguments which we have touched on in his discussion of the Turing Test to the main characteristics of his view of (3.1) the relation between the mental and the physical, (3.2) his epistemology, and (3.3) his theory of interpretation. After this I am going to discuss how those understandings come together in his unified theory of intellectual and linguistic competence and relate to Davidson's Test. Wrapping up with a discussion how Davidson's view might fit into the more recent discussion between connectionism and the classical representational and computational theory of mind.

### 3.1 The Meaning of Programs and Anomalous Monism

We will investigate Davidson's view of the relation between the mental and the physical following his essay 'Representation and Interpretation' ([1990a] 2004). Davidson starts out with the question "how much like us [a computer] must [...] be, and in what ways, to qualify as having thoughts" ([1990a] 2004, p. 88). As we have seen in subsection 2.3, he holds that a computer "cannot have thoughts unless it can learn and has learned from causal interactions with the world" ([1990a] 2004, p. 88). So "what must be added to computers [...] to insure that they are capable of thought" ([1990a] 2004, p. 89) — that they know *what* they are talking about — is "a very rich conceptual system" ([1990a] 2004, p. 90). Davidson's main question for the rest of the essay now becomes *how* we can judge that a computer has such a rich conceptual system. Davidson starts by ruling out that knowledge of the operation of the system is sufficient for this judgement:

"It may seem obvious that if an artificial object thinks and acts enough like a person, someone who knows how the object was designed and built would be able to describe and explain the mental states and actions of the object. But this does not follow, for there is no reason to suppose that there are definitional or nomological connections between the concepts used by the designer and the psychological concepts to be described and explained. This should be clear if we imagine that the builder has simply copied, molecule by molecule, some real person. [...] So one sort of 'complete' understanding does not necessarily imply another." (Davidson [1990a] 2004, p. 90)

The problem is one of *mental representation*. The "representation of an object or a fact [...] in the program cannot automatically be interpreted as a representation of that object or fact for the device" ([1990a] 2004, p. 91) — in its 'mind'. Davidson argues that indeed the representations differ. If an object or event falls under a certain mental concept there must not be an equivalent concept in the program or in a physical description. He gives the following analogy: "Suppose, following folk advice, I am attempting to go to sleep by counting sheep. Every now and then, at random, a goat slips into the file. In my drowsy state I find I cannot remember the classificatory words 'sheep' and 'goat'. Nevertheless I have no trouble identifying each animal: there is animal number one, animal number two, and so on. In my necessarily finite list, I can specify the class of sheep and the class of goats: the sheep are animals 1, 2, 4, 5, 6, 7, 8, and 12; the goats are animals 3, 9, 10, and 11. But these classifications

are no help if I want to frame interesting laws or hypotheses that go beyond the observed cases, for example, that goats have horns. I can pick out any particular sheep or goat in my animal numbering system, but I cannot, through conceptual poverty, tell the sheep from the goats generally. So it may be with the mental and the physical. Each mental event, taken singly, may have (must have, if I am right) a physical description, but the mental classifications may elude the physical vocabularies.” ([1990a] 2004, p. 92)

The program of a computer is purely syntactic, it can give explanations of its concepts in terms of other concepts it has, but it cannot specify anything about semantics, about references of its concepts, about relations of its concepts to objects or events in the world. “There is the language in which each animal can be picked out, but which lacks the concepts needed for classifying the animals as sheep or goats; similarly, syntax can provide a unique description of each true sentence, [...] but it can’t classify sentences as true or false. [...] If we knew no more than the program, we would have no reason to say [...] that any aspect of or event in the device represented anything outside the device” ([1990a] 2004, p. 91). Physics work in the same way, when it appeals to strict laws. Such strict laws are only “drawing on concepts from the same conceptual domain and upon which there is no improving in point of precision and comprehensiveness. [...] Physical theory promises to provide a comprehensive closed system guaranteed to yield a standardized, unique description of every physical event” ([1970] 2001, p. 223–224). Just like the program a complete physics is precise and contains no *ceteris paribus* (other things being equal) clauses in its explanations and it does not appeal to any concepts outside of its realm<sup>16</sup>. Therefore its semantics, its relation to the world is not given in the physical model or theory itself. Only our interpretation of physics gives its terms a semantic. This is realized through experimental practices in which we compare the predictions of the model to our observations using set standards<sup>17</sup>. We can try to describe the procedures of experiments as precisely as possible but not in the vocabulary of physics, they will always contain instructions in mental terms.

This is different for theories involving mental events; where “an event is mental if and only if it has a mental description [...] that involves terms like believing, intending, desiring, knowing, perceiving, remembering, and so on], the distinguishing feature of the mental is [...] what Brentano called intentionality” ([1970] 2001, p. 211). Mental properties are dispositional predicates which are causal in nature: “if Cain killed Abel, he must have done something that caused Abel’s death, and if he killed Abel intentionally, he must have been caused to act by a desire for Abel’s death. Beliefs, desires, and intentions are themselves causal dispositions” ([1990a] 2004, p. 95). Such causal dispositions do not feature in physics: “it is a sign of progress in a science that it rids itself of causal concepts. The dissolving of some salt is explained, up to a point, by saying that salt is soluble and this salt was placed in water; but one could predict the dissolving on the basis of far more general knowledge if one knew the mechanism, what it is about the constitution of the salt that accounts for its dissolving. When the mechanism is known, the explanation will not call on the causal concept of solubility” ([1990a] 2004, p. 96). This is also related to the semantical nature of our mind. In “specifying the contents of a memory, [we are forced] to refer to causes normally outside the person, [...] giving] a causal explanation of a belief [...] and] a semantic interpretation of it” ([1990a] 2004, p. 96). In this way mental representation is externalistic, it refers to terms outside of its descriptive scope.

So we have seen that physics and programs provide comprehensive theories that are purely syntactical and thus refer to no external con-

cepts which allow them to be expressed in strict laws which eliminate causal dispositions. We have seen that mental concepts on the other hand are semantic in an essential way and are expressed as causal dispositions and that such causal dispositions need to be eliminated to describe something in a complete physics as strict laws. Thus we can conclude that mental predicates can not occur in strict laws (See McLaughlin 2013, p. 417–423). This is what Davidson calls the Anomalism of the Mental: “there are no strict [...] laws on the basis of which mental events can be predicted and explained” (Davidson [1970] 2001, p. 208).<sup>18</sup>

But (other than McLaughlin 2013 seems to suggest) this is not the final argument Davidson has in mind, because externalism, *ceteris paribus* laws, and causal dispositions “could also be said to distinguish the concepts of many of the special sciences such as biology, [...] there must be something more basic or foundational” (Davidson [1991b] 2001, p. 217). Otherwise we would assume that “knowledge of the program of a device that successfully mimics the workings of the mind, will explain the mechanisms that support or constitute thought” (Davidson [1990a] 2004, p. 96–97) just like physics can explain the causal disposition of solubility.

“It is only when we can see a creature (or ‘object’) as largely rational by our own lights that we can intelligibly ascribe thoughts to it at all [...]. This means that when anyone [...] ascribes thoughts to others, he necessarily employs his own norms in making the ascriptions. There is no way he can check whether his norms are shared by someone else without first assuming that [the other is intelligible and] in large part [...] his norms] are [...] shared. [...] If the subject under study is to remain thought and intelligence, a normative methodology cannot be avoided.” (Davidson [1990a] 2004, p. 97–98)

“We depend on our linguistic interactions with others to yield agreement on [...] the sort of structures in nature that allow us to represent those structures in [...] numbers. We cannot in the same way agree on the structure of the sentences or thoughts [...], for the attempt [...] sends us back to the very process of interpretation [...]. It is here, I suggest, that we come to the ultimate springs of the difference between understanding minds and understanding the world as physical. A community of minds is the basis of knowledge; it provides the measure of all things. It makes no sense to question the adequacy of this measure, or to seek a more ultimate standard.” (Davidson [1991b] 2001, p. 218)

The interpretation of mental concepts is essentially normative, it requires us to assume a basic shared understanding of the world without which we can’t interpret others. We must assume others are largely coherent in their thoughts and corresponding to the same stimuli in the world as we are. This irreducibly normative character does not appear in physics or programs but it is this character which allows us to render any intelligible picture of the world. Physics and programs can’t give an account of their semantics, only our evaluation of their predictive powers through experiments puts them in a relation to the world. The concept of causality allows us to bridge the gap between these explanatory schemes. Its causal looseness is necessary to overcome the distance between the normative and formal explanations (see [1990a] 2004, p. 98). The ability to put things in relation to the world only emerges through our communication with others (we will discuss this more in the next subsection 3.2). The mental is irreducible in so far as it constitutes the basis of our very ability to make the world intelligible through such communication about a shared world.

We have seen that Davidson argues at length for the anomalism of the mental — mental events don't feature in strict laws and are thus irreducible to the physical. But we have also already seen that Davidson agrees with Turing that brains might act similarly to computers and that it should be generally possible to make thinking computers. In fact, Davidson rejects any cartesian mind body dualism. "The mental and the physical share *ontologies*, but not [...] classificatory concepts" ([1990a] 2004, p. 92).<sup>19</sup> Davidson labels this combination of ontological monism and agreement to the principle of the anomalism of the mental 'Anomalous Monism' (see "Mental Events" [1970] 2001). This view rests on the type-token-distinction: "each particular mental event [token] is identical with a particular physical event [token, but...] there are no strict psychophysical laws [...] and therefore] no mental event type is a physical event type" (Lepore and Ludwig 2009, p. 17).<sup>20</sup> So while the behavior can be completely explained on a physical level — since mental event tokens are identified with physical event tokens which are governed by strict laws — the physical explanation or the understanding of the program does not offer any insights for a mental explanation — since the physical event tokens don't fall under any mental event type. "Knowing the program is enough to explain why the device produces the marks or sounds or pictures it does given an input described in similarly abstract terms. This knowledge does not touch on questions [...] of reference to the outside world" (Davidson [1990a] 2004, p. 93). To make sense of the explanation by relating it to the world we need a mental scheme.

"Anomalous monism resembles materialism in its claim that all events are physical, but rejects the thesis [...] that mental phenomena can be given purely physical explanations. [...] Although the position I describe denies there are psychophysical laws, it is consistent with the view that mental characteristics are [...] dependent, or supervenient, on physical characteristics. Such supervenience might be taken to mean that [...] an object cannot alter in some mental respect without altering in some physical respect. Dependence or supervenience of this kind does not entail reducibility through law or definition" (Davidson [1970] 2001, p. 214)

Davidson's idea of the supervenience<sup>21</sup> and irreducibility of the mental to the physical has been highly discussed<sup>22</sup> and was and is very influential. As mentioned before it also provides the basis for the claim that "the only way to tell if an artificial device [...] can think], is to attempt to interpret the behavior of the device [...]because] interpretation involves the use of normative concepts [...] which] have no role in the understanding of a syntactically specified program" ([1990a] 2004, p. 99). So a behavioral test like Turing's and Davidson's Test is a necessary condition to judge anything as intelligent, be it a person, an animal, or an object.

### 3.2 Triangulating Knowledge of a Shared World

Similarly to the previous section we are going to examine Davidson's epistemology from a discussion that starts with another question that arises from the Turing test. This time, it is not the question *how* we can judge the computer, but *what* it is that would make a computer think. Davidson takes this route in 'What Thought Requires' ([2001] 2004):

"So we need to ask what would turn calculation, in the sense in which [...] a computer can calculate, into thought? [...] Some non-human animals can learn a great deal, but they do not learn that something is true. [...] Only a creature with propositional attitudes is equipped to fit a new concept into a complex scheme in which concepts have logical and other relations to one another. [...] Some degree of holism goes with having concepts.

The fundamental distinction [lies] between a mindless disposition to respond differentially to the members of a class of stimuli, and a disposition to respond to those items as members of that class. [...] A creature that cannot entertain the thought that it may be wrong has no concepts, no thoughts. To this extent, the possibility of thought depends on the idea of objective truth, of there being a way things are which is not up to us. [...] What is needed is something that can provide a standard against which an individual can check his or her reactions, and only other individuals can do this. Adding a second person helps [...] to pick out the relevant cause and to explain error]. It narrows down the relevant cause to the nearest cause common to two agents who are triangulating the cause by jointly observing an object and each other's reactions. The two observers don't share neural firings or incoming photons; the nearest thing they share is the object prompting both to react in ways the other can note. [...] Triangulation also creates the space needed for error, not by deciding what is true in any particular case, but by making objectivity dependent on intersubjectivity." (Davidson [2001] 2004, p. 138–137 & 141–143)

Davidson answers the question what distinguishes our ability to think from the ability to simply do very complex tasks like computers with the fact that we have concepts that occur in a holistic way. It is the concept of error that is the key to thought. The concept or error allows us to distinguish between pure conditioned responses (which can't be wrong, just like the calculations of a computer) and the right application of a concept. This requires semantics for a language which can distinguish between what we think and what is in the world. The core of the concept of semantics is the truth predicate; the idea that our thoughts and sentences can be true or false. This is also called the idea of objective truth — of a truth that lies beyond one's own judgement. This is the quintessential question of natural language philosophy raised by Wittgenstein: "what is the difference between using words correctly and merely thinking that one is using them correctly?" (Davidson [1994] 2005, p. 116)

The answer in Wittgenstein's famous argument against the possibility of a private language (1953, 256–282) is that language necessarily has a social aspect without which semantics can't be explained. Davidson relates to this. His essential argument is that to make the distinction between true and false believes a social background is required. Only because we hear different opinions we can call our own concepts into doubt.<sup>23</sup> The social norm that solves the problem is not given by conventions or shared practices, but by correlating reactions to a shared environment through a communication process.

"As Wittgenstein says, by yourself you can't tell the difference between the situations seeming the same and being the same. [...] If you and I can each correlate the other's responses with the occurrence of a shared stimulus, however, [...] it provides each of us with a ground for distinguishing the cases in which it fails. Failed natural inductions can now be taken as revealing a difference between getting it right and getting it wrong [...]. A grasp [...] of the distinction between thinking something is so and its being so,

depends on the norm that can be provided only by interpersonal communication" (Davidson [1994] 2005, p. 125)

Moreover<sup>24</sup>, Davidson argues that not even a physicalist account of thought can achieve that since it can't solve the problem of the inscrutability of reference (Quine [1960] 2013, ch. 2) — the unclarity what the relevant causal origin of a concept is (the object, the light, the stimulated neurons, ...). Davidson solves the problem by locating "the relevant cause of a speaker's thought [...] at the intersection of the two lines that can be drawn [...] from] the speaker [...] and] the interpreter [...] to] the object or event in the world they are responding to and communicating about" (Bernecker 2013, p. 445). Davidson fittingly calls this process of finding the intersection in the world between two points of view triangulation.

Since triangulation identifies the relevant cause of a thought in the distal stimulus in the world rather than in a proximal stimulus in our head (see Davidson [1990] 2005), it is a type of externalism. Different from Putnam's (1975) physical externalism and Burge's (1979) social externalism which try to identify the semantics of certain types of mental content and which Davidson rejects (Davidson:1987ko, [1994] 2005), his externalism is not concerned with the contents of any specific thoughts, but claims that interaction in a shared world is a necessary condition for having thoughts at all. This externalism has therefore been labeled "transcendental externalism" (Bridges 2006, p. 291).

### THREE VARIETIES OF KNOWLEDGE

This triangulation process is Davidson's epistemology. It can be captured as a relationship between three types of knowledge: (subjective) knowledge of the contents of one's own mind, (intersubjective) knowledge of the contents of other minds, and (objective) knowledge of the world. All three types of knowledge differ: Other than the knowledge of my own mind, knowledge of the world and other minds "depends on the functioning of my sense organs, and this causal dependence on the senses makes my beliefs [...] open to a sort of uncertainty" (Davidson [1991b] 2001, p. 205). My knowledge of the world is often "simply caused directly by the events and objects around me. But my knowledge [...] of other minds is never immediate in this sense" ([1991b] 2001, p. 205) it is only accessible through interpretation of behavior. The "three varieties of knowledge are concerned with aspects of the same reality; where they differ is in the mode of access to reality [...], but how can] the same world [be] known to us in three such different ways" ([1991b] 2001, p. 205 & 208)? This question lies at the heart of what Davidson identifies as the three main questions of epistemology: How can the mind have (1) knowledge about the natural world and (2) the content of another mind; and how can one (3) know the content of one's own mind without resort to evidence? ([1991b] 2001, p. 208)

Question 1 and 2 stem from the cartesian idea that knowledge of my own mind has conceptual priority over the others which leads to skepticism. Question 3 is raised by the positivists that answer it by reducing all knowledge to knowledge of world. But both accounts fail, since the three types are irreducible and none has conceptual priority. In Davidson's answer of triangulation all three questions are answered simultaneously. The content of propositional knowledge is directly caused by the objects and events the world is made up of. We come by this content through triangulating between our stimulus and that of another person caused by the same object or event in our shared world. We have access to other people's stimuli only

through interpretation of their (often verbal) reactions to those stimuli. Therefore all three types of knowledge are interdependent and none has priority above the others. My propositional knowledge (3) is dependent on interpretation of reactions to stimuli from a shared world. Knowledge about someone else's propositions (2) must be gauged by my own propositional knowledge. And knowledge about the world (1) is only possible through interpersonal communication.<sup>25</sup> The three form a triangle of interdependence. This leads Davidson to a denial of cartesian skepticism and positivist behavioral reductionism.

"Until a base line has been established by communication with someone else, there is no point in saying one's own thoughts or words have a propositional content. [...] The triangulation which is essential to thought requires that those in communication recognize that they occupy positions in a shared world. So knowledge of other minds and knowledge of the world are mutually dependent[...]

Attributing thoughts to others is a matter of matching the verbal and other behavior of others to our own propositions [...]. Knowledge of our own minds and knowledge of the minds of others are thus mutually dependent. [...]

The nature of interpretation guarantees both that a large number of our simplest perceptual beliefs are true, and that the nature of these beliefs is known to others. [...] Any particular belief [...] about the world around us may be false. What cannot be the case is that our general picture of the world and our place in it is mistaken, for it is this picture which informs the rest of our beliefs and makes them intelligible [...]." (Davidson [1991b] 2001, p. 213)

Let us turn to the last conclusion Davidson draws here. Interpretation is under-determined: to know whether a speaker holds a sentence true or not, we need to know what he means by it, which requires us to triangulate which in turn only works if we can determine whether he holds the sentence true. To break this circle we have to appeal to a norm of communication. Davidson claims that only an interpersonal standard of consistency and correspondence (sometimes also called Principles of Charity) can achieve this. "The Principle of Coherence prompts the interpreter to discover a degree of logical consistency in the thought of the speaker; the Principle of Correspondence prompts the interpreter to take the speaker to be responding to the same features of the world that the interpreter would be responding to under similar circumstances" ([1991b] 2001, p. 211). The former ensures that a holistic interpretation is possible and the latter allows us to create the bridge of triangulation between beliefs and references. Those principles are essential for interpretation and therefore also for knowledge. They make it "constitutive of what it is to be a speaker that one is mostly right about the external world, one's own thoughts, and what one's words mean. The importance of this conclusion can hardly be overemphasized. If it is right, then we have a transcendental guarantee of knowledge of our own minds, the minds of others [...], and the external world. We secure this without having to explain how it is that we justify our beliefs on the basis of evidence, for knowledge in each of these domains emerges as a fundamental condition on having the capacity to speak and think at all." (Lepore and Ludwig 2009, p. 20)

Understanding the importance and function of triangulation makes it clear why communication outside of a shared world does not allow us to say much about the nature of understanding of a speaker. Only through observing the interdependent three fold relationship

between interpreter, speaker, *and the world* we can tell anything about the semantics of a speaker. Davidson's emphasis in his test on the importance of observing the computer interact in a shared world in order to judge its intelligence, is rooted in this conviction about epistemology.

### 3.3 From Prior to Passing Theories of Interpretation

Now we finally come to the heart piece of Davidson's philosophy: his theory of interpretation. I will first give a rough sketch of the main idea of the theory and then discuss its anti-conventionalism and the formalizations Davidson proposes and how those seemingly contradictory accounts might be reconciled.

Davidson's interest lies not in "the (empirical) question how we actually go about understanding a speaker [...but in] the (philosophical) question what is necessary and sufficient for such understanding" (Davidson [1994] 2005, p. 111). He does not aim to give a full account of how we interpret others in everyday practice nor "to provide useful hints to real linguists, or to criticize their methods [...or] yield an insight into [...] first-language acquisition" (Davidson [1995a] 2004, p. 131). He is concerned with the essence of our linguistic abilities that allows us to interpret propositional content and gets to the core of what our mind can do. I adopt the term *linguistic competence*<sup>26</sup> as opposed to linguistic *practice* to refer to this. A theory of linguistic competence is a model that can give an account how interpretation of propositional contents of others can be described under some idealized conditions. In his early writing Davidson refers to this using the term 'Radical Interpretation' ([1973b] 2001), which originates from Quine's account of radical translation ([1960] 2013) — a *hypothetical* situation in which an interpreter figures out how to understand the previously completely unknown 'language' of an unknown speaker (similar *but idealized* to the task of a linguists learning a native's language). This account of interpretation without any previous knowledge of how to proceed is what gets to the core of linguistic competence according to Davidson (it is important to note that this is not about *acquisition* but about *description* of linguistic competence).

The core of Davidson's theory is that the process of interpreting utterances can be described by a theory that yields an interpretation for each utterance of the speaker that can be expressed in the interpreter's language. To understand someone we need to 'translate' his utterances into our own vocabulary. It becomes clear that some 'translation' is necessary even if we speak the same 'language' when we recognize that no two people assign exactly the same interpretations to all their words, so no two people will simply share the understanding of all their utterance, because they share a 'language'. What happens in the process of understanding each other is described by Davidson as the convergence of two speakers prior theories of interpretation into a passing theory of interpretation. The prior theory of interpretation of the interpreter is what he comes equipped with to interpret this particular speaker and for the speaker it is what informs his way of talking to this specific interpreter. The passing theory is what they converge on in their understanding that allows them to interpret the utterance in the same way.

"I have distinguished [...] the prior theory from [...] the passing theory. For the hearer, the prior theory expresses how he is prepared in advance to interpret an utterance of the speaker, while the passing theory is how he *does* interpret the utterance. For the speaker, the prior theory is what he believes the interpreter's prior

theory to be, while his passing theory is the theory he *intends* the interpreter to use." (Davidson [1986] 2005, p. 101)

Sometimes the speaker's and interpreter's prior theories of interpretation for a sentence might already converge, but very often we need to do some mental work to arrive at a shared passing theory (you might take basically any sentence from this paper as an example that you do not understand instantly). Therefore *radical* interpretation of sorts is quite common in everyday communication. What does it mean to have the same passing theory for an utterance? It means that we assign the same meaning to an utterance. For Davidson this is true if the same truth conditions would be assigned to the utterance by speaker and interpreter (see [1965] 2001, more on this later). Davidson's main question is not how we get a prior theory of interpretation in the first place (he assumes that we have some linguistic ability already), but how we converge on a passing theory of interpretation. "Linguistic ability is the ability to converge on a passing theory" ([1986] 2005, p. 107) of interpretation. Therefore, to account for linguistic competence we must have a meta-theory that describes how to transition a prior into a passing theory. This is where Davidson's holism comes in. We need to appeal to propositional attitudes like believe, preference, and intention to interpret a speaker. Davidson's description of such a meta-theory has two parts: In the first step we determine the syntax and logical form based on simple preferences and patterns. In the second step we determine the meanings of simple expressions through triangulation in a shared environment and then from there find the meanings of complex terms. It is important not to misunderstand Davidson as claiming that speakers or interpreters use theories to literally translate each other: "the point is not that [a] speaker or hearer *has* a theory, but that they speak and understand *in accord with* a theory — a theory that is needed *only* when we want to *describe* their abilities" ([1994] 2005, p. 113, emphasis added).

Because of the distributed nature of Davidson's work, the reconstruction of his theory of interpretation is a bit cumbersome and controversial. Two of the most prominent traits of his theory that appear in his early and late essays are the rejection of meanings as entities sui generis or as expressions in a language of thought and an accompanying anti-conventionalism that holds that no rules known to speaker and hearer in advance can be enough to converge on a passing theory. As we have seen in the previous section about triangulation — Davidson holds truth and communication as the elementary concepts which make meaning and conventions circumstantial. Nevertheless he thinks that Interpretation can be treated with formalisms and is as serious of a science as any. In the following I will discuss how these views fit together.

#### THE SOCIAL ASPECT OF LANGUAGE BEYOND CONVENTIONS

In his essays 'Communication and Convention' ([1984] 2001) Davidson examines the relation between meaning and use of language for communication. Since Austin ([1956] 1979) pointed out that speech acts have not only a *descriptive* character but also a *performative* character that *changes* the state of the world, a main topic of the Philosophy of Language has been to describe how those two aspects are related. A common way to explain that relation is to appeal to conventions. Following Austin ([1955] 1962), Davidson distinguishes between three aspects of an utterance: *locution* (literal meaning), *illocution* (intended meaning), and *perlocution* (ulterior purpose) ([1984] 2001, p. 272). He argues against the proposal that the relation between literal meaning and propositional attitudes (like

purpose and intention) can be established by appeal to conventions. "Grammatical mood and illocutionary [and perlocutionary] force, no matter how closely related, cannot be related simply by convention. [...] Convention can [not] link what our words mean — their literal semantic properties, including truth — and our purposes [and intention] in using them" ([1984] 2001, p. 271).<sup>27</sup>

We should not be surprised by Davidson's anti-conventionalism since we have already seen that he rejects the idea that knowledge of other minds is reducible. Trying to explain illocution and perlocution (which depend on propositional attitudes) with conventions in terms of the grammatical mood and meaning is trying to do just that. This becomes pretty clear if we think about using conventions to help a computer communicate without having any other means for interaction than linguistic ones — which is what Davidson denies. He goes even further in denying "that the meaning of a word is conventional, that is, that it is a convention that we assign the meaning we do to individual words and sentences when they are uttered or written" ([1984] 2001, p. 276). A regularity in assignment of meaning can only hold if speaker and hearer have the same assignment of meaning to the utterance prior to the speech act. However the prior and passing theory of interpretation is often not the same. Therefore no convention for the assignment of meaning exists, but the assignment develops in the act of communication.

"What must be shared for communication to succeed is the passing theory. For the passing theory is the one the interpreter actually uses to interpret an utterance, and it is the theory the speaker intends the interpreter to use. [...] But the passing theory cannot in general correspond to an interpreter's linguistic competence. [...] Every deviation from ordinary usage, as long as it is agreed on for the moment [...], is in the passing theory as a feature of what the words mean on that occasion. Such meanings [...] are what I have called first meanings. [...] Of course things previously learned were essential to arriving at the passing theory, but what was learned could not have been the passing theory." (Davidson [1986] 2005, p. 102-103)

"As the speaker speaks his piece the interpreter alters his theory, entering hypotheses about new names, altering the interpretation of familiar predicates, and revising past interpretations of particular utterances in the light of new evidence. [...] We can not have all this knowledge previous to the utterance, because [...] a speaker may provide us with information relevant to interpreting an utterance in the course of making the utterance." (Davidson [1986] 2005, p. 101)

Any prior shared practice, convention, or meaning is not essential for understanding, because there is "no reason [...] why speakers who understand each other ever need to speak [...] as anyone else speaks" ([1994] 2005, p. 115).<sup>28</sup> The only notion of a shared practice that Davidson accepts is the assumption "that two speakers couldn't understand each other if each couldn't (pretty well) say in his way what the other says in his" ([1994] 2005, p. 115-116). In his essay 'The Social Aspect of Language' ([1994] 2005), Davidson clarifies what he takes to be the underlying reason for this rejection. Any linguistic "conformity is contingent on the desire to be understood. [...] The threat of failure to communicate is the reason for conforming" ([1994] 2005, p. 118) to any linguistic rules, not the other way round. "If we want to be understood, all we need to worry about is how our actual audience will take our words. [...] I would not speak the words I do if I thought they would not be understood" ([1994] 2005, p. 118 & 123). Without basic linguistic abilities no shared practice can

be established nothing is more basic than successful communication, meaning emerges from it, not the other way around.

"The intention to be taken to mean what one wants to be taken to mean is [...] clearly the only aim that is common to all verbal behavior [...]. Success in communicating propositional contents [...] is what we need to understand before we ask about the nature of meaning or of language [...]. Meaning [...] gets its life from those situations in which someone intends that his words will be understood in a certain way, and they are. [...] Where understanding matches intent we can, if we please, speak of 'the' meaning; but it is understanding that gives life to meaning, not the other way around." (Davidson [1994] 2005, p. 121-122)

Knowing a language, i.e. the shared and previously mastered "ability to operate in accord with a precise and specifiable set of syntactic and semantic rules [...], is neither necessary nor sufficient for successful linguistic communication" ([1994] 2005, p. 110).<sup>29</sup> Davidson agrees with Wittgenstein that the problem how to distinguish between using words correctly and merely thinking one is using them correctly is crucial and that it can only be solved by appeal to a social environment. But he disagrees that the norm for getting it right is given by conventions or shared practices. As we have seen in subsection 3.2 it is solved by Triangulation which does not require any shared practice but does depend on interactions in a shared world and in turn does not depend on the ability to interpret but gives rise to it.<sup>30</sup> That an interpreter can understand a speaker if he talks about a salt shaker on the table is not because the word 'salt shaker' has a clear definition or a clearly defined conventional usage that we can imagine as something living in a platonic world of ideas. On the contrary, the meaning of the utterance 'salt shaker' only arises if the interpreter understands what the speaker is referring to and they can use the word to successfully interact in the world. The word could be replaced by anything else like 'salt thing' or just 'thing' and would still have the same meaning if the interpreters and speakers communication would succeed. Meaning is not a condition but a result of successful communication.

Even though Davidson rejects that convention has any essential role in explaining interpretation, it is "a practical crutch to interpretation, a crutch we cannot in practice afford to do without — but a crutch which [...] we can in the end throw away, and could in theory have done without from the start" ([1984] 2001, p. 279). Especially for determining syntax and grammatical mood convention are often used, since grammar is way more socially stable than semantics. "If we can apply our general method of interpretation to a speaker [...] it is] because we can treat his structure-forming devices as we treat ours [...], fixing] the logical form of his sentences, and [determining] the parts of speech." ([1984] 2001, p. 279)

## A RECURSIVE AND EMPIRICAL THEORY OF INTERPRETATION

Davidson sets the framework for his investigation of linguistic competence through two criteria that a theory of interpretation must fulfill. First, he holds that no a priori considerations can suffice to describe the mechanisms of interpretation of a natural language. We need to study interpretation empirically. For an empirical study of interpretation "it must be possible to give a constructive account of the meaning of the sentences in the language. Such an account [is what Davidson] calls a theory of meaning for the language" (Davidson [1965] 2001, p. 3). Second, the crucial requirement on such a

theory is that it is finite — “a learnable language has a finite number of semantical primitives” ([1965] 2001, p. 9).<sup>31</sup> If we “regard the meaning of each sentence as a function of a finite number of features of the sentence” ([1965] 2001, p. 8) we can explain how we can recursively determine the meaning of infinitely many sentences by only knowing the meaning of a finite number of expressions. So Davidson’s wants to describe linguistic competence through an empirical and recursive theory of interpretation.

The problem of meaning goes back to Frege ([1892] 1993). While we can interpret words as representations for objects in the world, it is unclear what sentences refer to. We need a proper theory of the meaning of sentences. But many existing theories of meaning do not fulfill Davidson’s requirements ([1965] 2001, p. 9–16). In ‘Truth and Meaning’ ([1967] 2001) Davidson argues that introducing meanings as entities *sui generis* does not help in the construction of a *recursive* theory of meaning. Following Quine (1953) Davidson also refutes that a linguists account of meanings based on analyzing the syntax and then using a dictionary or lexicon to look up meanings of words can work, since it obviously fails for believe sentences and is circular. A theory is needed that can determine the meanings of words holistically from their use in sentences. Davidson finds his solution in Tarski’s truth theory (1944), taking the notion of truth as basic.<sup>32</sup> Davidson believes that the notion of meaning can be replaced by simply giving the truth conditions for a sentence (relative to the speech event to capture indexical expressions). Instead of looking for what a sentence means we look for interpretations of sentences which can be given by other sentences of explicative nature which we can formally generate with Tarski’s theory.<sup>33</sup>

It is important to not misunderstand Davidson as claiming that the meaning of a sentence is its truth condition.<sup>34</sup> As we have seen before, correct interpretation — successful communication — is the essential aspect of language. Thus Davidson talks about ‘Radical Interpretation’ ([1973b] 2001) and holds that an interpretation of a sentence is only justified holistically as part of an empirical theory that optimally fits evidence about the speakers meanings, believes, and desires ([1973b] 2001, p. 133–140). In that sense Davidson fully embraces Austin’s account of the performative aspect of language. Interpretation is as much about meaning as it is about beliefs and desires, all three can only be determined simultaneously (Rescorla 2013, p. 474). But to give a formal account of the prior and passing theory of an interpreter a formal device is needed to capture the content of utterances. Since truth conditions describe under which circumstances in the world something is true, truth conditions are well fit to give a description of the content of an utterance in a way that is relatable to observations about the world. And such observations about the world are what ultimately determines the correctness of the theory empirically. Because the believes of the speaker are of key importance, but not accessible by direct evidence. The interpreter must compare his believes to the interpreted truth conditions and whatever else he has as evidence about the speakers believes and desires. This can only work if the interpreter applies a principle of charity. This principle assumes that the speaker largely shares the interpreters logic and understanding of the world. “Interpretation must ascribe a background of true beliefs to the speaker [...] and depict the speaker as largely conforming to the interpreter’s own rational norms, including norms given by logic, probability theory, and decision theory” (Rescorla 2013, p. 477).<sup>35</sup> A theory of linguistic competence must be a theory about rationality and intellectual competence as well.

Inspired by his his early collaboration with Patrick Suppes and JJC McKinsey in Decision Theory at Stanford (Lepore and Ludwig 2009,

p. 2), Davidson developed his formal account of Radical Interpretation further into ‘A Unified Theory of Thought, Meaning, and Action’ ([1980] 2004).

“[The] three structures, of logic, decision theory, and formal semantics, have the characteristics of serious theories in science: they can be precisely [...] stated, and, given empirical interpretation and input, they entail endless testable results. Furthermore, logic, semantics, and decision theory can be combined into a single unified theory of thought, decision, and language, as I have shown. This is to be expected. Decision theory extracts from simple choices subjective scales for probabilities, i.e., degrees to which sentences are held to be true, and for values or the extent to which various states of affairs are held to be desirable. Radical interpretation [...], extracts truth conditions, that is, meanings and beliefs, from simple expressions of assent and dissent. Formal semantics has logic built in, so to speak, and so does decision theory in the version of Richard Jeffrey. Uniting the theories depends on finding an appropriate empirical concept, and one such concept is the relation between an agent, the circumstances of utterance, and two sentences, one of which the agent would rather have true than the other.” (Davidson [2001] 2004, p. 146–147)

Davidson’s model consists of two steps. In the first step the syntax (logical structure and grammatical mood) is determined, in the second step the semantics (relations to world experiences) are determined. Davidson believes that both can be formalized using Tarski’s truth theory and Bayesian decision theory which can “codify belief through subjective probability and desire through utility [...] and] show how to extract a semi-unique theory of mental states (probabilities and utilities) from a relatively observable evidentiary base (preferences). The theorems thereby illuminate how one can read a rational pattern into observed behavior” (Rescorla 2013, p. 475).

Davidson gives a fairly clear outline of the first step; explaining how logical connectives can be extracted from evidence about preferring the truth of one sentence over another and than determining logical quantifiers from patterns of those preferences.<sup>36</sup> After which the grammatical moods can be determined with a Tarski style theory. So the first step yields a complete analysis of the logical form of a sentence for a rational agent.

The “only remaining task is to interpret primitive singular terms and predicates. Davidson’s remarks regarding this crucial task are not as systematic as one might desire” (Rescorla 2013, p. 475). This second step is where the described triangulation process comes into play to give an interpretation to the logical form of the sentence that is expressed in truth conditions about the world. But how exactly a Bayesian theory can be applied to this remains mostly unclear from Davidson’s writing.<sup>37</sup>

## LINGUISTIC COMPETENCE, PRACTICE, AND META-THEORY

If we reflect on the just given description of Davidson’s empirical and fairly formalized two step unified theory of linguistic and intellectual competence, it seems unclear at first how it could be reconciled with remarks like this:

“There are no rules for arriving at passing theories, no rules in any strict sense, as opposed to rough maxims and methodological generalities. A passing theory really is like a theory at least in this, that it is derived by wit, luck, and wisdom from a private vocabulary and grammar, knowledge of the ways people get their

point across, and rules of thumb for figuring out what deviations from the dictionary are most likely. There is no more chance of regularising, or teaching, this process than there is of regularising or teaching the process of creating new theories to cope with new data in any field — for that is what this process involves.” (Davidson [1986] 2005, p. 107)

His unified theory doesn't seem very dependent on luck and wit, but more like an applicable formal tool which yields interpretations. We should first remember that his theory is not aiming to describe our linguistic practice, but give a formalized account to help shed light on the essence of linguistic competence.

“The approach to the problems of meaning, belief, and desire that I have outlined is not, I am sure it is clear, meant to throw any direct light on how in real life we come to understand each other [...]. I have been engaged in a conceptual exercise aimed at revealing the dependencies among our basic propositional attitudes at a level fundamental enough to avoid the assumption that we can come to grasp them [...] one at a time. My way of performing this exercise has been to show how it is in principle possible to arrive at all of them at once.” (Davidson [1980] 2004, p. 166)

This comment from the paper where he first introduces his formal account of the unified theory should lead our attention to another point. While the unified theory contains a formalized way to derive the syntax it does so only based on empirical evidence and holistically. It does not allow to derive even the syntax based only on the utterances, but it requires knowledge about preferring sentences true. But furthermore the second step in which the semantics are determined is much more holistic and requires knowledge of world as well. So it is not a reductive theory in any way.

Davidson understands the unified theory as a sort of proof or at least plausibility check for the process of triangulation and radical interpretation. He is convinced that philosophy should be a serious science that aims to produce empirical theories. The unified theory is a sketch of a formal account of triangulation that can be tested empirically. Maybe we should interpret it as his attempt to establish a real alternative to a positivistic and behavioristic account of the essence of linguistic and intellectual competence — in the spirit of the idea that without an alternative no scientific paradigms can shift (Kuhn 1962). In his essay ‘Could There Be a Science of Rationality?’ ([1995a] 2004) he examines whether his account of rationality can be considered a serious science. He identifies “three features of the Unified Theory [...] that have been thought to remove it from the domain of serious science [...]: its assumptions of holism and of externalism, and its normative properties” ([1995a] 2004, p. 129). He concludes that those properties are not sufficient to prove the theory as unscientific, but that they make the theory irreducible. As we have seen in subsection 3.1 this has to do with a missing independent standard to ground the evaluation of the theory in. “We cannot check up on the objective correctness of our own norms by checking with others, since to do this would be to make basic use of our own norms once more” ([1995a] 2004, p. 134). This might be the key to understand his remarks that we cannot give a fully explanatory account of the meta theory that describes how to transition from prior to passing theories; for no account will give a reductive analysis. The application of the unified theory relies on the principle of charity, it relies on using our already existing understanding of the world. Insofar it is not formal tool but depends on our wit and luck of understanding the world. It stays somewhat circular but it is a broad and informative circle that helps to grasp our linguistic and intellectual competence.

### 3.4 Davidson's Test and the Theory of Intellectual and Linguistic Competence

- Interpreting Davidson's Test as an expression of his theory of linguistic and intellectual competence

#### WHAT IS MEANING: INTERPRETATIONISM AND EXTERNALISM

<- - Davidson rejects that we need meaning, it ought to be replaced with translation - His interpretationism seems to not sit well with his externalism, but the test provides a solution - Meaning is not primary to communication, but the other way round. (see next section) - Davidson rejects Putnam's physical externalism with Anomalous monism, thoughts are in the head (see Bernecker) - Davidson rejects Burge's social externalism with his anti-conventionalism - The Underlying Ontology and Metaphysics - Ontology of only objects and events (E3.14, E5.4) - Directness of knowledge of world - no intermediary entities (5.3, E5.4, E3.10) - no mentales, language only through communication, meaning given causally by objects events (E5.9) - Truth as elementary non-reducible concept (E5.2, E2.1) ->

“an externalist view [...] is a straightforward consequence of his interpretationist stance, together with the view that evidence is [given by] distal, not proximal stimuli” (Hahn 2003, p. 35)

Interpretationism: there is only a gradual, not a type difference between the interpretation of a believe and holding a believe. In that sense meanings are expressions, not things in the world. (Schwitzgebel [2006] 2015)

How does such interpretationism sit with triangular externalism? We might have to “treat them as independently necessary and jointly sufficient conditions for having thoughts. Having thoughts requires not only that one is interpreted, but also that one learned words in the right sort of way” (Bernecker 2013, p. 447)

38

“The only other way we can bring in causal history is as simply an extra condition not on interpretability, but on having thoughts. Being interpreted is not sufficient for having thoughts; the subject's language must also have been learned in the right way. But this, congenial as it might be to someone like Burge perhaps, seems deeply antithetical to Davidson's insistence that the third-person interpretive stance is all there is to the attribution of meanings and thoughts: that it is constitutive of them.” (Hahn 2003, p. 45)

Other than Hahn suggests this interpretation is explicitly the one Davidson takes in ‘Turing's Test’.

“The intention to be taken to mean what one wants to be taken to mean is [...] clearly the only aim that is common to all verbal behavior [...]. Success in communicating propositional contents [...] is what we need to understand before we ask about the nature of meaning or of language [...]. Meaning [...] gets its life from those situations in which someone intends that his words will be understood in a certain way, and they are. [...] Where understanding matches intent we can, if we please, speak of ‘the’ meaning; but it is understanding that gives life to meaning, not the other way around.” (Davidson [1994] 2005, p. 121-122)



## OBJECTIONS AGAINST DAVIDSON'S TEST

- behaviorism
- mechanism
- scope of the test
- Objection based on Internal Operation Arguments
- Blockheads
- The Chinese Room (Internal Operation)
- French: Associative Priming and Rating Games
- abstract terms and wirkungsgeschichte

### Ability to cope with Abstract Terms

It remains unclear how well Davidson's theory might be able to cope with abstract terms in the language.

He throws very little light on how understanding of abstract terms might work.

There are good reasons to doubt that communication about abstract terms can work in the same way as other communication. Something like the Wirkungsgeschichte becomes way more plausible, if we look at how culturally loaded language is

See also Bertram's critique: Gadamer's idea (Bertram 2012)

Davidson to Gadamer, (Davidson [1997a] 2005)

## ANTI-REPRESENTATIONALISM AND CONNECTIONISM

- status of program not interpretable
- non-deterministic in some way
- externalistic because of data dependence
- Davidson as a connectionist, rejects classicist Computational Theory of Mind, agrees with computability of mind
- Also externalist, anomalous monism solves problem of external dependence of mind, even though brain is in the skull

"A neural network is a collection of interconnected nodes. Nodes fall into three categories: input nodes, output nodes, and hidden nodes (which mediate between input and output nodes). Nodes have activation values, given by real numbers. One node can bear a weighted connection to another node, also given by a real number. Activations of input nodes are determined exogenously: these are the inputs to computation." (Rescorla 2015)

## Conclusion

Davidson's theory investigates "what it is about propositional thoughts [...] that makes them intelligible to others. This is a question about the nature of thought and meaning which cannot be answered by discovering neural mechanisms, studying the evolution of the brain, or finding evidence that explains the incredible ease and rapidity with which we come to have a first language." (Davidson [1995a] 2004, p. 133)

### Davidson's Theory and Machine Learning

- is Davidson's argument of understanding a program really justified

While Davidson is skeptical about the possibility of explaining our linguistic performance. He believes that there is a pragmatic and empirical theory about the essence of linguistic competence.

## Implications for Artificial Intelligence and Computer Linguistics

Davidson tries to devise a theory that can model human intellectual abilities in an empiric theory that has the power to explain them. (As opposed to only a simulation of the brain by a machine)

While it currently seems more probable that Turing's Test and also Davidson's Test will be passed by a computer that leverages trained large scale deep neural networks aka. something into the direction of simulating the brain. Davidson's consideration still provide important input into what is needed to achieve human level of natural language abilities. That is the ability to somehow interact with the world and have some sensory impression of it and have a triangular model of determining world interaction.

- Possibility of AI and the importance of Language
- Davidson sees artificial intelligence as possible.
- He agrees that language is the key test for intelligence.
- Davidson and Turing agree that the task is not to explain human linguistic practice, but create an empirical theory that models the essence of linguistic competence. This is a very important philosophical distinction that should be understood very well in order to build language processing systems and not get distracted with the wrong problems.
- Important boundary conditions for linguistic competence
- Linguistic and intellectual competence are inherently social/intersubjective
- Linguistic and intellectual competence require interaction with the shared world (empiric)
- Denial of conventionalism and statistical linguistics
- The social and empiric traits of linguistic and intellectual competence can not be modeled by conventionalism, if this convention requires a previously shared common practice. Interpretation emerges in the instance of communication and cannot be predetermined in anyway.
- Abilities of linguistic interaction can not be clearly separated from abilities for physical interaction. Both are necessary for intelligence.

## Notes

1 — It is remarkable that a contemporary philosopher of Davidson's calibre only started publishing his relevant papers in his 40s after he had been a professor for over a decade. The publication dates of Davidson's papers ought to be interpreted with care though. His friend and student Ernie Lepore recounts: "[Davidson was] adventurous and daring [...] from early on right up until his death. [...] Donald was without even the slightest speck of careerism from the very start. He traveled the world giving papers in exotic places and often handed them over to local journals upon request. [...] Many of the papers were written somewhat contemporaneously. [...] Some were] given much earlier than they were published and relatively around the same time. [...] I possess a mimeographed copy of his quotation paper dated from the early 60's, long before its 1979 publication." (Lepore 2003).

Davidson's life and career are actually quite inspiring, *fascinating*, and unconventional. Davidson started out in the History of Philosophy and ventured into psychology and economy during his studies and as a professor. Certainly it is no coincidence that he was cited by continental and analytical philosophers and even computer scientists alike. To learn more about Davidson's life I recommend to read Lepore's 'Interview with Donald Davidson' ([1999] 2004) and Davidson's own 'Intellectual Autobiography' (1999a)

2 — I think it is pretty clear that Davidson ([1990b] 2004, p. 78) misinterprets Turing when he suggests that Turing wants the computer to play the literal imitation game and pretend to be a *woman* and that the interrogator's task would be to decide on the gender. This becomes pretty clear when Turing says: "If the man were to try and pretend to be the machine [...]" (1950, p. 435) and from all his examples that are focused on how the machine can imitate being a *human* not a woman. (See also Copeland 2000, p. 526.)

3 — Turing's idea is that a baseline is established by the traditional imitation game, in which a man tries to imitate a woman and the interrogator has to decide on the gender. While this might not be the best way to establish a baseline, it means that Davidson's critique that "Turing does not say what he would make of a computer that was consistently chosen over the [...human] to be the [...human]" ([1990b] 2004, p. 78), is ill-conceived, as the computer will unambiguously fall over or under the baseline within a margin of error.

4 — The term 'intelligence' is adopted here purely in reference to Artificial Intelligence. Turing mostly uses the term 'thinking' and Davidson ([1995a] 2004) mostly uses the term 'rationality' to refer to the concept that this thesis is concerned with. I will introduce and explicate the term intellectual competence for this concept in subsection 3.4.

5 — See Muehlhauser (2013) for a brief overview of operational definitions of AGI, including the Coffe-Brewing and College Test. From a philosophical perspective all those seem rather random and dubious — certainly much further away from clearly capturing necessary conditions of human-like intellectual competences.

6 — Shieber (1994) criticizes the Loebner Prize for its inappropriateness to award advances in natural-language-processing techniques instead of engineering tricks oriented to the exigencies of the restricted task like parrying and insertion of random typing errors. The setup of the scoring system alone shows how pointless it is to even judge current systems by a direct Turing Test. He argues that a subjective award modeled after the Nobel Prizes would make significantly more sense.

7 — This is also known as the 'Church-Turing-Deutsch Principle' and represents an extension of the well-known Church-Turing Thesis to artificial intelligence.

8 — Steven Wolfram writes: "Computers (and especially linear algebra in GPUs) got fast enough that [...] it became practical to train neural networks with millions of neurons, on millions of examples. [...]his suddenly brought large-scale practical applications within reach. [...] I don't think it's a coincidence that this happened right when the number of artificial neurons being used came within striking distance of the number of neurons in relevant parts of our brains.

[...]f we're trying to achieve 'human-like' image identification [...] then this defines a certain scale of problem, which, it appears, can be solved with a 'human-scale' neural network." (2015)

9 — For historical reasons, the terminology here is easily confusing: The term 'Empiricism' ought to be treated very carefully as distinct from empirical. Davidson's theory is empirical, meaning observable evidence plays an important role in it, but it disagrees with the dogmas of Logical Empiricism (sometimes also called Logical Positivism), namely the analytic-synthetic distinction and reductionism (Quine 1951) and the dualism of scheme and content (Davidson [1973a] 2001). Davidson replaces the concepts of meaning and proximal stimuli by a holistic and empirical theory which treats language and knowledge as elementary intersubjective. (Lepore and Ludwig 2009, p. 22-23)

10 — Davidson moved from his interest in the historic aspect of philosophy to analytical philosophy, because Quine had convinced him "that it was possible to be serious about getting things right in philosophy" and to take philosophy "as serious as science" (Davidson and Lepore [1999] 2004, p. 239). But through the recession of Logical Empiricism it had become clear that this progress could not simply lie in treating philosophy as the logic of mathematics and reasoning. His experiences in studying business, working with J.C.C. McKinsey on decision theory, and studying Tarski, gave Davidson "an appreciation for what it's like to have a serious theory" (Davidson and Lepore [1999] 2004, p. 253) and convinced him that the progress of ordinary language philosophy could be found in looking for a pragmatic empirical theory.

11 — Davidson argues that: "If we have the semantics of a language right, the objects we assign to the expressions of the language must exist" (Davidson [1993] 2005, p. 40). Since being able to communicate successfully means getting semantics pretty right. This will in turn also mean, that if we can communicate successfully, we must get it pretty right what things are in the world and in which relation they stand to each other. And this certainly seems to be a sufficient condition for intelligence. Davidson's argument for his claim that semantics is a method for metaphysics ([1993] 2005) is based on the fact that we must get things mostly right about the world, if we are intelligible as rational beings at all.

12 — Davidson's ontology takes only objects and events as basic entities ([1991b] 2001). He rejects the need for any propositions of proximal stimuli as entities in the mind ([1990] 2005). The content of propositional knowledge is directly caused by the distal objects and events the world is made up of.

13 — In fact Descartes somewhat preconceived Turing's test in his 'Discourse on the Method'. He claims: machines "could never use words or other signs arranged in such a manner as is competent to us in order to declare our thoughts to others" (Descartes [1637] 1993, Part V).

14 — Davidson gives the following example to illustrate his point: "Alitalia Flight 19 leaves Turin for London on Tuesdays at 8:30 in the morning. We can learn this by consulting a computer; but does the computer know what we learn by consulting it? The answer is that it does not because it does not know what a flight is, where Turin is, or even that Tuesday is a day of the week." ([1990a] 2004, p. 89)

15 — To illustrate this point, let me give an example from first semester physics: Take momentum in classical physics — the fact that we single out the meaning of this word by attributing it as a part of the empirical theory of classical mechanics already shows that it gains its relevance from its place in this theory. It gains its relevance from its place in this empirical theory that serves to describe the movement of mass under natural and artificially applied forces. We can't debate the definition of momentum ( $\vec{p}$ ), not because it is defined in very clear mathematical terms as the time derivative of the product of mass and location ( $\vec{p} = \partial_t(m \cdot \vec{r})$ ), but because changing its definition would destroy some of the predictive power of classical mechanics. If we would for example only take the simplified definition of momentum as the product of mass and velocity (which assumes the mass doesn't change over time ( $\partial_t m = 0 \Rightarrow \vec{p} = m \cdot \partial_t \vec{r} = m \cdot \vec{v}$ ) — which is still very clearly defined mathematically — we couldn't explain how a rocket flies that accelerates by throwing away mass ( $\partial_t m \neq 0$ ); the theory of classical mechanics would lose some of its power to fit empirical evidence. This is closely related to Popper's (1935) requirement of falsifiability for a scientific theory. The definitions are only good, as long as changing them would mean that the theory as a whole can't fit some empirical evidence anymore. In this sense their meaning and relevance is determined holistically.

16 — This is assuming a complete physics of the world, which the practical science of physics strives to reach. In the science of physics of course there are still many issues like the inability of quantum mechanics to define in physical terms what an observation is that leads to the collapse of the wave function. But the fact that this is seen as a huge deficiency of the theory by most physicists illustrates this description of a complete physics — independent of whether it is achievable or not.

17 — Units are the bridges between the variables in the physical model and our observable reality. Their definitions give a semantics to the model. But their definitions are externalistic. This becomes clearest if we look at the definition of the basic SI unit 'kilogram' which is currently still defined by a weight lying around in a vault in Paris — not by the number of atoms of it, but by the object that we can use in comparative experiments.

18 — In his section 2.(8) about the 'Argument from Informality of Behaviour' Turing discusses a somewhat related argument that there is no "definite set of rules of conduct" (1950, p. 452) which govern the behavior of a person and that therefore a deterministic computer could not simulate a person's behavior. As Turing points out this argument (which he is not deriving directly from any other literature, but establishing himself in this crude way) fails because it is beside the point and has no evidence for its claim. But unlike Oppy and Dowe ([2003] 2011, p. 21, a fairly poor entry in the Stanford Encyclopedia of Philosophy) suggest the argument could easily be made strong if we suppose a mind-body dualism, causal influence of the mind on the body, and one of the various arguments like the one given here that there are no strict rules governing mental processes.

19 — Davidson argues that 'Spinoza's Causal Theory of the Affects' ([1999] 2005) might be interpreted as a form of anomalous monism.

20 — Another famous argument against reductionism by Fodor (1974) is also related to the fact that a special science type has multiple realizations through physical tokens, but those will fail to establish any physical type predicate at least in most cases.

21 — The modern use of the term supervenience actually originates in 'Mental Events' (Davidson [1970] 2001, p. 214). "This is the first explicit statement of a psychophysical supervenience thesis in the literature" (McLaughlin 2013, p. 427, see also p. 438).

22 — Davidson's 'Mental Events' ([1970] 2001) has been a very influential and highly discussed paper. For an introduction to some main discussions around Anomalous Monism about mental and physical taxonomies and event individuation, strict laws and completed physics, Ramsification, strong and weak supervenience, and type epiphenomenalism and the cause law principle see McLaughlin (2013).

23 — Some have argued that autistic speakers show that having the idea of objective truth and interpreting others is not essential for the ability to think and speak. There is empirical evidence that "autistic speakers allegedly fail false belief tests, which are designed to test their ability to attribute false beliefs to others [...and] actual cases of autistic speakers not capable of attributing beliefs to others have reportedly been found." Since those autists are clearly capable of language and thought, it is argued that the social aspect of language is unnecessary. It remains very unclear how convincing the evidence from such tests is and what alternative picture of language without a social environment and the idea of objective truth could possibly look like. (For a brief introduction to the discussion and overview of the literature, see Verheggen 2013, p. 466 & 469.)

24 — It has been suggested that two main arguments for triangulation can be extracted from Davidson's writing. The first one, "the argument from error", is the one trying to solve the origin of the idea of objective truth which I have just sketched; the second one, the "argument from object-directedness", tries to solve the problem of the inscrutability of reference which I am about to outline now (Bridges 2006, p. 292). My aim here is not to reconstruct these arguments, but to plausibilize Davidson's triangulation idea. For an introduction to some of the main problems of Davidson's arguments see Bernecker (2013) and Bridges (2006) and for some defenses see Verheggen (2013).

25 — Some have argued that Davidson's account of thought in triangulation is circular since language presupposes thought (for example see Bernecker 2013, Bridges (2006)). Davidson himself is "under no illusion that [he] can provide anything like an analysis; [...] for a non-circular answer would tell us how to account for intensionality in non-extensional terms" (Davidson [1997b] 2005, p. 139) which he sees as impossible as we have seen in subsection 3.1. Nevertheless the account is instructive. "It is obviously not just the narrow vicious circle that says that there is language and thought only when there is language and thought. It is rather the rich and complex circle that encompasses language, thought (each of which depending on there being fixed meanings), possession of the concept of objective truth, and linguistic triangulation. I should say it is an instructive circle, each part of which is such that, if it is absent, the whole circle breaks down. [...] Therefore] it seems to me that the story is compelling and that within the intensional realm, progress has been made" (Verheggen 2013, p. 468). So the critique is somewhat beside the point. The relevance of any theory does not follow from its ability to reduce terms but from its ability to account for empirical data (otherwise String Theory would be accepted). Davidson's account is relevant because it fits well with many of our experiences how we gain knowledge and is therefore an attractive alternative to cartesian idealism or positivist reductionism.

26 — The term linguistic competence was introduced by Chomsky who makes a “fundamental distinction between competence (the speaker-hearer’s knowledge of his language) and performance (the actual use of language in concrete situation)” (Chomsky 2014, p. 4). While Davidson criticizes the idea of linguistic competence as knowledge of a language, he adopts the term for the description of the interpreters essential competence that enables communication in ‘A Nice Derangement of Epitaphs’ ([1986] 2005).

27 — The core of the arguments Davidson establishes against a conventional connection between illocutions or perlocutions and meanings takes assertion as an example and shows that there can not be a convention (public sign) that signals that an utterance is sincere (claims the truth): “In making an assertion, the asserter [...] must intend that [his] intention [to make an assertion] be recognized by his audience. [...] So it is natural to think it would be useful if there were a convention, as a convenience in making our assertive intentions clear. But Frege was surely right when he said, ‘There is no word or sign in language whose function is simply to assert something.’ [...] Imagine this: [an] actor is acting a scene in which there is supposed to be a fire. [...] And now a real fire breaks out, and the actor tries vainly to warn the real audience. [...] Convention cannot connect what may always be secret — the intention to say what is true — with what must be public — making an assertion” ([1984] 2001, p. 269–270). But without a convention that connects truth values with intentions or purposes of an utterances there is no hope in explaining meaning (a function from utterances to truth values) by those properties of utterances and conventions or the other way round. This is not to be confused with “the (correct) thesis that every utterance of an imperative [or declaration of believe] labels itself (truly or falsely) an order [or assertion]” ([1984] 2001, p. 275).

28 — Davidson’s prime example for instances where the ability to interpret only arises in the moment of the utterance are Malapropisms. “Malapropisms introduce expressions not covered by prior learning, [...] they fall into a different category, one that may include such things as our ability to perceive a well-formed sentence when the actual utterance was incomplete or grammatically garbled, our ability to interpret words we have never heard before, to correct slips of the tongue, or to cope with new idiolects.” (Davidson [1986] 2005, p. 94–95)

29 — Davidson has become famous for claiming “that *there is no such thing as a language* [...]. We must give up the idea of a clearly defined shared structure which language-users acquire and then apply to cases.” ([1986] 2005, p. 107, emphasis added). Our discussion should illuminate what he means by that. Of course he does not argue that the word ‘language’ is meaningless. He says “speakers share a language if and only if they tend to use the same words to mean the same thing[. Which leaves...] defining a language as the philosophically rather unimportant task of grouping idiolects.” ([1994] 2005, p. 111)

30 — Relation to Plato and Gadamer (Davidson [1997a] 2005, Bertram (2012)). Relation to Derrida (Bertram 2002, Wheeler III (1986))

31 — Davidson argues that a language must be finite, because otherwise it would take us an infinite amount of time to learn a language: “Suppose that a language lacks this feature; then no matter how many sentences a would-be speaker learns to produce and understand,

there will remain others whose meanings are not given by the rules already mastered. It is natural to say such a language is unlearnable” ([1965] 2001, p. 8). Even though this argument might have some issues, since the number of expressions in a language is only infinite in theory, but in practice we ever only use a finite number of words. However, it is clear that theories that treat for example Quotations as irreducible singular terms go completely against our intuition, so Davidson certainly has a point with his critique of such theories.

32 — Davidson later realized that he reversed Tarski’s approach: “One thing that only gradually dawned on me was that while Tarski intended to analyse the concept of truth by appealing (in Convention T) to the concept of meaning (in the guise of sameness of meaning, or translation), I have the reverse in mind. I considered truth to be the central primitive concept, and hoped, by detailing truth’s structure, to get at meaning” (Davidson 2001b, p. xvi). This is of course related to Davidson taking successful communication as basic for meaning, where communication depends on triangulating in a shared world, allowing one to determine what is objectively true. So we can see that these later ideas are already moored in his early works.

33 — Davidson’s idea that truth conditions (T-Sentences) are all that is needed to explain interpretation, and any more fundamental notion of meaning is misconceived, has also led to what is known as Davidson’s Program: The attempt to give an account of the interpretation of utterances purely with the means of first order logic on which Tarski’s truth theory is based. Davidson has made important contributions to such analyses himself. Especially his analysis of action sentences with first order logic has gained widespread use even by computational linguists. “Action sentences present a problem for semantics because of their capacity to take an endless variety of adverbial phrase, which themselves can be endlessly complex. [...] The event analysis that Davidson introduced treats action verbs as introducing an implicit existential quantifier over events, and the adverbs as contributing predicates of it.” (Lepore and Ludwig 2009, p. 16) The sentence “Brutus stabbed Caesar violently with a knife” for example can analyzed as:  $\exists e : e = \text{‘event of stabbing by Brutus of Caesar’} : \text{UsingKnife}(e) \wedge \text{Violent}(e)$ .

34 — Davidson clarifies the following: “I do not think I have ever conflated the (empirical) question how we actually go about understanding a speaker with the (philosophical) question what is necessary and sufficient for such understanding. I have focused on the latter question [...]. If I ask how someone interpreted an utterance of the sentence ‘Snow is white’, and am told that she interpreted it as meaning that snow is white (or as being true if and only if snow is white), my question was not, as the answer shows, what other words the hearer might have substituted for the sentence ‘Snow is white.’ I am asking how the person understood the utterance of those words. Of course I must use words to say how she understood those words, since I must use words to say anything, but my words are not offered as the interpretation; they merely help describe it. [...] I agree with Michael that ‘one who ...understands a sentence need not be able to say how he understands it. He does not have to be able to say it even to himself’ [Dummett (1986), p. 464 ...]. It would obviously have been absurd of me to have claimed [...] that whenever we understand a speaker we translate his words into our own. Translation is no part of the transaction between speaker and hearer that I call interpretation. Where translation of a sort may be involved is in the description the philosopher gives in his language of what the hearer makes of the speaker’s utterances. [...] Similarly] the

point is not that speaker or hearer has a theory, but that they speak and understand in accord with a theory — a theory that is needed only when we want to describe their abilities and performance” (Davidson [1994] 2005, p. 111–113).

35 — Davidson has sometimes been accused of behaviorism for his account of radical interpretation (see for example Chomsky 1992). This is clearly mistaken, as Davidson sees behavior purely as evidence and holds that intentional terms are irreducible to behavioristic vocabulary. He says “Chomsky has accused me [...] of supposing that all we know about language must be based on behavioristic evidence. I would certainly deny the accusation [...]. But I do share with Quine the conviction that our understanding of what speakers mean by what they say is partly based, directly or indirectly, on what we can learn or pick up from perceiving what they do. No matter how much grammar we come equipped with from the cradle, we must learn what the words of any particular language mean [...]; we must pick up our first language from those who already speak it. (The behaviorism I speak of is not, incidentally, reductive in nature: I do not expect any basic intentional predicates to be defined in non-intentional terms. The point simply concerns evidence)” (Davidson [1995a] 2004, p. 132). On the contrary Davidson is a strong supporter of the importance of rational norms for the explanation of cognitive phenomena (Rescorla 2013).

36 — Davidson uses Bayesian decision theory in the form developed by Richard Jeffrey which allows “to extract subjective probabilities and values from preferences that propositions [or in Davidson’s modification sentences] be true” (Davidson [1980] 2004, p. 160). This allows to construct a formal theory that assigns values of probability that an agent believes in a sentences and values of desirability of the truth of a sentence for an agent. The clue is that these values do not have to be measured directly; “once the truth-functional connectives have been identified, Jeffrey has shown how to fix [...] the subjective desirabilities and probabilities of all sentences” ([1980] 2004, p. 161). Davidson formally shows how we can determine which logical connective has the meaning of the Sheffer stroke (‘not both’) based purely on the knowledge whether a speaker prefers the truth of one sentence over that of another ([1980] 2004, p. 161–163). Since all logical operators can be reduced to the Sheffer stroke this means we have already determined all logical operators. From there “the relative desirabilities of all sentences” are determined [p. 164] and from the patterns of those sentences the logical quantifiers can be identified. Uncovering “logical form in general, that is, to learn how sentences are made up of predicates, singular terms, quantifiers, variables, and the like” ([1980] 2004, p. 165) only based on observations whether an agent prefers the truth of one sentence over that of another.

However, as “ingenious and sophisticated [as this proof is, ...] the three-place relation between the speaker and any two sentences [preferring the truth of one over the other...] is not an attitude we are familiar with [...] without knowing what the speaker thinks or means” (Glock 2003, p. 192–193). What is gained until this problem is solved? More than Glock (2003) suggests, I would say. If we are inclined to believe that the relation could be observed, for example in the kind of betting scenarios Davidson describes ([1980] 2004, p. 161–162), we have a guarantee of identifying the logical structure in the unknown language of a coherent speaker and have a foundation of a theory of interpretation in the three basic concepts of believe, desire, and meaning (in a Tarskian style). Unfortunately it seems that little work has been done to find ways of specifying how the preference relationship can be observed or in testing the empirical strength

of Davidson’s theory, so this discussion remains largely speculative. It is also important to note that — as we have seen in subsection 3.1 — Davidson does not think that the whole theory can be captured in strict laws, so some degree of looseness is to be expected.

37 — It is a point of critique that “Davidson develops [...his theory] primarily through armchair reconstruction of interpretive practice, rather than through detailed study of empirical psychology” and cognitive science (Rescorla 2013, p. 484). This leads to a disconnect between his aim for a serious scientific treatment of intellectual and linguistic competence (see Davidson [1995a] 2004) and the standard of his theory. It is especially said because “over the past few decades, cognitive science has embraced the Bayesian paradigm [...that] embodies a broadly Davidsonian [...] entanglement of normative evaluation and psychological description” (Rescorla 2013, p. 484). If Davidson had more explicitly engaged with cognitive science his theory might have become clearer and cognitive scientists might have paid more explicit attention to the relevance of his philosophy for their research.

38 — The discussion about these issues often takes Davidson’s Swampman example ([1987] 2001) as a starting point. I will not go any deeper into the example since Davidson himself later regretted using it since it caused more confusion than anything else: “As with Swampman, I regret these sorties into science fiction [...] the case can be made without [such sorties], and better” (1999b, p. 192).

## Acknowledgements

Many thanks to ... for reading this paper ...

### COLOPHON

This paper was written with Gingko in Pandoc Markdown and typeset using XeTeX in Gotham Narrow. My research process is inspired by Luhmann’s Zettelkasten methodology, using Evernote. If you are interested, you can find an article about my research and writing process on my blog (<http://MrLoh.se/2015/10/research-and-scientific-writing-process>).

## References

- AISB. 2015. “Loebner Prize.” <http://www.aisb.org.uk/events/loebner-prize>.
- Austin, J. L. (1955) 1962. *How to Do Things with Words*. Oxford University Press.
- Austin, J. L. (1956) 1979. “Performative Utterances.” In *Philosophical Papers*, p. 233–252. Oxford University Press.
- Bernecker, Sven. 2013. “Triangular Externalism.” Edited by Kirk Lepore and Ernie amd Ludwig. Wiley-Blackwell.
- Bertram, Georg W. 2002. “Übergangsholismus — Holismus, Veränderung Und Kontinuität in Den Sprachphilosophien von Davidson

- Und Derrida." *Zeitschrift Für Philosophische Forschung* 56 (3). Vittorio Klostermann: p. 388–413.
- Bertram, Georg W. 2011. *Sprachphilosophie Zur Einführung*. Junius-Verlag.
- Bertram, Georg W. 2012. "Antwort und Zugang." In *Hermeneutik Und Die Grenzen Der Sprache*, edited by Ulrich Arnsward, Jens Kertscher, and Louise Röska-Hardy. Heidelberg.
- Bickle, John. (1998) 2013. "Multiple Realizability." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta, Spring 2013. <http://plato.stanford.edu/archives/spr2013/entries/multiple-realizability/>.
- Block, Ned. 1981. "Psychologism and Behaviorism." *The Philosophical Review* 90 (1). Duke University Press: p. 5–43. <http://www.jstor.org/stable/2184371>.
- Bridges, Jason. 2006. "Davidson's Transcendental Externalism." *Philosophy and Phenomenological Research* 73 (2): p. 290–315.
- Burge, Tyler. 1979. "Individualism and the Mental." *Midwest Studies In Philosophy* 4 (1). Blackwell Publishing Ltd: p. 73–121. doi:10.1111/j.1475-4975.1979.tb00374.x.
- Butterfill, Stephen A, and Ian A Apperly. 2013. "How to Construct a Minimal Theory of Mind." *Mind & Language* 28 (5). Blackwell Publishing Ltd: p. 606–637.
- Cavell, Marcia. 2004. "Introduction." In *Problems of Rationality: Philosophical Essays Volume 4*, p. xiii–xix. Oxford University Press.
- Chomsky, Noam. 1992. "Language and Interpretation, Philosophical Reflections and Empirical Inquiry." In *Inference, Explanation, and Other Frustrations: Essays in the Philosophy of Science*, edited by John Earman, 14:p. 99–128. University of California Press.
- Chomsky, Noam. 2014. *Aspects of the Theory of Syntax*. MIT press.
- Copeland, B Jack. 1999. "A lecture and two radio broadcasts on machine intelligence by Alan Turing." In *Machine Intelligence 15, Intelligent Agents [St. Catherine's College, Oxford, July 1995]*, p. 445–476. Oxford University.
- Copeland, B Jack. 2000. "The Turing test." *Minds And Machines* 10 (4): p. 519–539.
- Copeland, B Jack, and Diane Proudfoot. 1999. "Alan Turing's Forgotten Ideas in Computer Science." *Scientific American*, p. 99–103.
- Davidson, Donald. 1999a. "Intellectual Autobiography." In *The Library of Living Philosophers: The Philosophy of Donald Davidson*, edited by Lewis E Hahn, p. 3–70. Open Court.
- Davidson, Donald. 1999b. "Reply to a.C. Genova." In *The Library of Living Philosophers: The Philosophy of Donald Davidson*, edited by Lewis E Hahn, p. 192–194. Open Court.
- Davidson, Donald. (1963) 2001. "Actions, Reasons, and Causes." *Journal of Philosophy* 60 (23): p. 685–700.
- Davidson, Donald. (1974a) 2001. "Belief and the Basis of Meaning." *Synthese* 27 (July-August): p. 309–323.
- Davidson, Donald. (1984) 2001. "Communication and Convention." *Synthese* 59 (1): p. 3–17.
- Davidson, Donald. (1991a) 2001. "Epistemology Externalized." *Dialectica* 45 (23): p. 191–202.
- Davidson, Donald. 2001a. *Essays on Actions and Events: Philosophical Essays Volume 1*. Oxford University Press.
- Davidson, Donald. 2001b. *Inquiries Into Truth And Interpretation: Philosophical Essays Volume 2*. Oxford University Press.
- Davidson, Donald. (1987) 2001. "Knowing One's Own Mind." In *Subjective, Intersubjective, Objective: Philosophical Essays Volume 3*, p. 15–38. Clarendon Press.
- Davidson, Donald. (1970) 2001. "Mental Events." In *Experience and Theory*, edited by L F Oster and J W S Wanson, p. 79–101. Humanities Press.
- Davidson, Donald. (1973a) 2001. "On the Very Idea of a Conceptual Scheme." *Proceedings and Addresses of the American Philosophical Association* 47 (n/a): p. 5–20.
- Davidson, Donald. (1973b) 2001. "Radical Interpretation." *Dialectica* 27 (1): p. 314–328.
- Davidson, Donald. (1974b) 2001. "Replies to David Lewis and W.V. Quine." *Synthese* 27 (3-4): p. 345–349.
- Davidson, Donald. (1976) 2001. "Reply to Foster." In *Truth and Meaning: Essays in Semantics*, edited by Gareth Evans and John Henry McDowell, p. 33–41. Clarendon Press.
- Davidson, Donald. 2001c. *Subjective, Intersubjective, Objective: Philosophical Essays Volume 3*. Clarendon Press.
- Davidson, Donald. (1988) 2001. "The Myth of the Subjective." In *Bewusstsein, Sprache Und Die Kunst*, edited by M K Rausz, p. 221–240. Notre Dame University Press.
- Davidson, Donald. (1992) 2001. "The Second Person." *Midwest Studies in Philosophy* 17 (1): p. 255–267.
- Davidson, Donald. (1965) 2001. "Theories of Meaning and Learnable Languages." In *Proceedings of the International Congress for Logic, Methodology, and Philosophy of Science*, edited by Yehoshua Bar-Hillel, p. 3–17. North-Holland.
- Davidson, Donald. (1991b) 2001. "Three Varieties of Knowledge." In *Royal Institute of Philosophy Supplement*, edited by A P Hillips Giffiths, p. 153–166. New York: Cambridge University Press.
- Davidson, Donald. (1967) 2001. "Truth and Meaning." *Synthese* 17 (3): p. 304–323.
- Davidson, Donald. (1978) 2001. "What Metaphors Mean." *Critical Inquiry*, no. 5: p. 31–47.
- Davidson, Donald. (1995a) 2004. "Could There Be a Science of Rationality?" *International Journal of Philosophical Studies* 3 (1): p. 1–16.
- Davidson, Donald. (1985) 2004. "Incoherence and Irrationality." *Dialectica* 39 (4): p. 345–354.
- Davidson, Donald. 2004. *Problems of Rationality: Philosophical Essays Volume 4*. Oxford University Press.
- Davidson, Donald. (1990a) 2004. "Representation and Interpretation." In *Modelling the Mind*, edited by K V Newton-Smith W H and Wilkes, p. 13–26. Oxford University Press.
- Davidson, Donald. (1995b) 2004. "The Problem of Objectivity."

*Tijdschrift Voor Filosofie* 57 (2): p. 203–220.

Davidson, Donald. (1980) 2004. "Toward a Unified Theory of Meaning and Action." *Grazer Philosophische Studien* 11: p. 1–12.

Davidson, Donald. (1990b) 2004. "Turing's Test." In *Problems of Rationality: Philosophical Essays Volume 4*. Oxford University Press, Originally published in *Modelling the Mind*, edited by K Said, Oxford University Press: 1990.

Davidson, Donald. (2001) 2004. "What Thought Requires." In *The Foundations of Cognitive Science*, edited by João Bragança, p. 121. Oxford: Clarendon Press.

Davidson, Donald. (1986) 2005. "A Nice Derangement of Epitaphs." In *Truth and Interpretation: Perspectives on the Philosophy of Donald Davidson*, edited by Ernest Lepore, p. 433–446. Blackwell.

Davidson, Donald. (1997a) 2005. "Gadamer and Plato's Philebus." In *The Philosophy of Hans-Georg Gadamer*, edited by L Hahn. Chicago.

Davidson, Donald. (1990) 2005. "Meaning, Truth, and Evidence."

Davidson, Donald. (1993) 2005. "Method and Metaphysics." *Deucalion*, no. 11: p. 239–248.

Davidson, Donald. (1997b) 2005. "Seeing Through Language." In *Truth, Language and History: Philosophical Essays Volume 5*, p. 127–142.

Davidson, Donald. (1999) 2005. "Spinoza's Causal Theory of the Affects." In, edited by Yirmiahu Yovel. Little Room Press.

Davidson, Donald. (1994) 2005. "The Social Aspect of Language." In *The Philosophy of Michael Dummett*, edited by Gianluigi McGuinness Brian and Oliveri, p. 1–16. Kluwer.

Davidson, Donald. 2005. *Truth, Language and History: Philosophical Essays Volume 5*. Oxford University Press.

Davidson, Donald, and Gilbert Harman. (1970) 2001. "Semantics of Natural Language." *Synthese* 22 (1–2): p. 1–2.

Davidson, Donald, and Ernest Lepore. (1999) 2004. "An Interview with Donald Davidson." In *Problems of Rationality: Philosophical Essays Volume 4*, p. 231–266. Oxford: Oxford University Press. [http://rucss.rutgers.edu/faculty/lepore/images/Davidson\\_interview.pdf](http://rucss.rutgers.edu/faculty/lepore/images/Davidson_interview.pdf).

Descartes, René. (1637) 1993. Project Gutenberg.

Dummett, Michael. 1986. "A Nice Derangement of Epitaphs: Some Comments on Davidson and Hacking." In *Truth and Interpretation: Perspectives on the Philosophy of Donald Davidson*, edited by Ernest Lepore, p. 459–476. Blackwell Publishers.

Fodor, Jerry A. 1974. "Special Sciences (or: The Disunity of Science as a Working Hypothesis)." *Synthese* 28 (2). Springer: p. 97–115.

Frege, Gottlob. (1892) 1993. "Über Sinn Und Bedeutung." In *Logische Untersuchungen*. Vandenhoeck & Ruprecht.

Garson, James. (1997) 2015. "Connectionism." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta, Spring 2015. <http://plato.stanford.edu/archives/spr2015/entries/connectionism/>.

Glock, Hans-Johann. 2003. *Quine and Davidson on Language, Thought and Reality*. Cambridge University Press.

Gómez-Torrente, Mario. (2006) 2015. "Alfred Tarski." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N Zalta. <http://plato.stanford.edu/archives/spr2015/entries/tarski/>.

Hahn, Martin. 2003. "When Swampmen Get Arthritis: 'Externalism' in Burge and Davidson." In *Reflections and Replies: Essays on the*

*Philosophy of Tyler Burge*, edited by Martin Hahn and Bjørn Ramberg, p. 29–58. MIT Press.

Horgan, Terence, and John Tienson. 1989. "Representations Without Rules." *Philosophical Topics* 17 (1): p. 147–174.

Jurafsky, Daniel, and James H Martin. 2015. *Speech and Language Processing*. 2nd ed. An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.

Kuhn, Thomas S. 1962. *The Structure of Scientific Revolutions*. The University of Chicago Press.

Lau, Joe, and Max Deutsch. (2002) 2014. "Externalism About Mental Content." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta, Summer 2014. <http://plato.stanford.edu/archives/sum2014/entries/content-externalism/>.

Lepore, Ernest. 2003. "Memorial Eulogy for Donald Davidson." [http://rucss.rutgers.edu/faculty/lepore/images/Davidson\\_memorial.pdf](http://rucss.rutgers.edu/faculty/lepore/images/Davidson_memorial.pdf).

Lepore, Ernest, and Kirk Ludwig. 2005. *Donald Davidson: Meaning, Truth, Language, and Reality*. Oxford University Press, USA.

Lepore, Ernest, and Kirk Ludwig. 2009. "Davidson." In *12 Modern Philosophers*, edited by Christopher Belshaw and Gary Kemp. Wiley-Blackwell.

Lepore, Ernie, and Kirk Ludwig. 2007. *Donald Davidson: Meaning, Truth, Language, and Reality*. Clarendon Press.

Malpas, Jeff. (1996) 2015. "Donald Davidson." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N Zalta. <http://plato.stanford.edu/archives/fall2015/entries/davidson/>.

McLaughlin, Brian P. 2013. "Anomalous Monism." In *A Companion to Donald Davidson*, edited by Kirk Lepore Ernie and Ludwig, p. 410–442. Wiley-Blackwell.

Millar, P Hartley. 1973. "On the Point of the Imitation Game." *Mind* 82 (328). Mind Assoc: p. 595–597.

Moor, James. 1976. "An Analysis of the Turing Test." *Philosophical Studies* 30 (4): p. 249.

Muehlhauser, Luke. 2013. "What is AGI?" <https://intelligence.org/2013/08/11/what-is-agi/>.

Oppy, Graham, and David Dowe. (2003) 2011. "The Turing Test." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N Zalta. <http://plato.stanford.edu/archives/spr2011/entries/turing-test/>.

Pitt, David. (2000) 2013. "Mental Representation." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta, Fall 2013. <http://plato.stanford.edu/archives/fall2013/entries/mental-representation/>.

Popper, Karl. 1935. *Logik Der Forschung. Zur Erkenntnistheorie Der Modernen Naturwissenschaft*. Wien: Springer-Verlag.

Putnam, Hilary. 1975. "The Meaning of 'Meaning'." *Minnesota Studies in the Philosophy of Science* 7. University of Minnesota Press: p. 131–193.

Quine, Willard Van Orman. 1951. "Two Dogmas of Empiricism." *The Philosophical Review* 60 (1). New York: Harper Torchbooks: p. 20–43.

Quine, Willard Van Orman. 1953. "The Problem of Meaning in Linguistics." *From a Logical Point of View*. Harvard University Press, p.

Quine, Willard Van Orman. (1960) 2013. *Word and Object*. 2nd ed. MIT Press.

Rescorla, Michael. 2013. "Rationality as a Constitutive Ideal." In *A Companion to Donald Davidson*, edited by Kirk Lepore and Ernie amd Ludwig, p. 472-488. Wiley-Blackwell.

Rescorla, Michael. 2015. "The Computational Theory of Mind." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta, Winter 2015. <http://plato.stanford.edu/archives/win2015/entries/computational-mind/>.

Schubert, Lenhart. (2014) 2015. "Computational Linguistics." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N Zalta, Summer 2015. <http://plato.stanford.edu/archives/spr2015/entries/computational-linguistics/>.

Schwitzgebel, Eric. (2006) 2015. "Belief." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N Zalta, Summer 2015. <http://plato.stanford.edu/archives/sum2015/entries/belief/>.

Searle, John R. 1980. "Minds, Brains, and Programs." *Behavioral and Brain Sciences* 3 (3). Cambridge University Press: p. 417-424.

Sher, Gila. 1999. "What Is Tarski's Theory of Truth?" *Topoi* 18 (2). Springer: p. 149-166.

Shieber, Stuart M. 1994. "Lessons from a restricted Turing test." *ArXiv Preprint Cmp-Lg/9404002*.

Speaks, Jeff. (2010) 2014. "Theories of Meaning." In *The Stanford Encyclopedia of Philosophy*, edited by Edward N Zalta. <http://plato.stanford.edu/archives/fall2014/entries/meaning/>.

Tarski, Alfred. 1944. "The Semantic Conception of Truth: And the Foundations of Semantics." *Philosophy and Phenomenological Research* 4 (3): p. 341-376.

Turing, Alan. 1950. "Computing Machinery and Intelligence." *Mind*

59 (October): p. 433-460.

Turing, Alan. (1948) 1992. "Intelligent Machinery." In *Collected Works of a. M. Turing: Mechanical Intelligence*, edited by D. C. Ince. Elsevier Science Publishers.

Turing, Alan. (1951) 1996. "Intelligent Machinery, a Heretical Theory." *Philosophia Mathematica* 4 (3): p. 256-260.

Turing, Alan. (1952) 1999. "Can Automatic Calculating Machines Be Said To Think?" In *A Lecture and Two Radio Broadcasts by Alan Turing*, edited by B J Copeland. In *Machine Intelligence* 15, p. 466-476, Oxford University Press.

Turing, Alan. (1951) 1999. "Can Digital Computers Think?" In *A Lecture and Two Radio Broadcasts by Alan Turing*, edited by B J Copeland. In *Machine Intelligence* 15, p. 462-465, Oxford University Press.

Verheggen, Claudine. 2013. "Triangulation." Edited by Kirk Lepore and Ernie amd Ludwig. Wiley-Blackwell.

Wheeler III, Samuel C. 1986. "Indeterminacy of French Interpretation: Derrida and Davidson." In *Truth and Interpretation: Perspectives on the Philosophy of Donald Davidson*, edited by Ernest LePore, p. 477-494. Blackwell Publishers.

Wikipedia. 2015. "Helen Keller." [https://en.wikipedia.org/wiki/Helen\\_Keller](https://en.wikipedia.org/wiki/Helen_Keller). (accessed: October 28, 2015).

Wittgenstein, Ludwig. 1953. *Philosophische Untersuchungen*.

Wolfram, Stephen. 2002. *A New Kind of Science*. Vol. 5. Wolfram Media.

Wolfram, Steven. 2015. "Wolfram Language Artificial Intelligence: The Image Identification Project." <http://blog.stephenwolfram.com/2015/05/wolfram-language-artificial-intelligence-the-image-identification-project>.